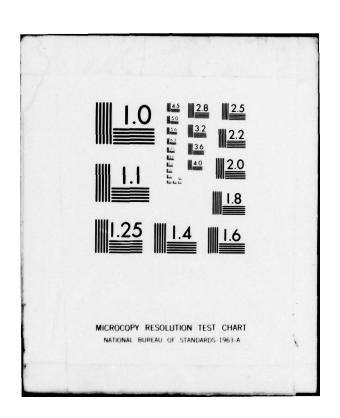
TEXAS INSTRUMENTS INC DALLAS
LIMITED VOCABULARY CONTINUOUS WORD RECOGNITION. (U)
AUG 79 R L DAVIS, B 6 SECREST, C J CATO F30
RADC-TR-79-204 AD-A074 206 F/6 17/2 F30602-77-C-0139 UNCLASSIFIED NL 1 of 2



RADC-TR-79-204
Final Technical Report
August 1979



LIMITED VOCABULARY CONTINUOUS WORD RECOGNITION

Texas Instruments, Inc.

Robert L. Davis Bruce G. Secrest Craig J. Cato DDC DCOCMIC SEP 25 1979 ULGGUTG

223

APPROVED FOR PUBLIC RELEASE; DISTRIBUTION UNLIMITED

MDA074206

UNG FILE WILL

ROME AIR DEVELOPMENT CENTER Air Force Systems Command Griffiss Air Force Base, New York 13441

79 39 24 001

This report has been reviewed by the RADC Information Office (OI) and is releasable to the National Technical Information Service (MTIS). At NTIS it will be releasable to the general public, including foreign nations.

RADC-TR-79-204 has been reviewed and is approved for publication.

Richard S. Vonusa

RICHARD S. VOMUSA Project Engineer

APPROVED:

hut h

ROSS H. ROGERS, Colonel, USAF Chief, Intelligence & Reconnaissance Division

FOR THE COMMANDER: John S. Huss JOHN P. HUSS Acting Chief, Plans Office

If your address has changed or if you wish to be removed from the RADC mailing list, or if the addressee is no longer employed by your organiza-tion, please notify RADC (IRAA), Griffise AFB NY 13441. This will assist us in maintaining a current mailing list.

Do not return this copy. Retain or destroy.

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

19 REPORT DOCUMENTATION PAGE	READ INSTRUCTIONS BEFORE COMPLETING FORM
RADC TR-79-204	NO. 3. RECIPIENT'S CATALOG NUMBER
LIMITED VOCABULARY CONTINUOUS WORD RECOGNITION.	5. TYPE OF REPORT & PERIOD COVERED Final Technical Report August 1977 — January 1979
	6. PERFORMING ORG. REPORT NUMBER
Robert L. /Davis / Bruce G./Secrest Craig J./Cato	8. CONTRACT OR GRANT NUMBER(s) F30602-77-C-0139
PERFORMING ORGANIZATION NAME AND ADDRESS Texas Instruments Inc. (13500 North Central Expressway Dallas TX 75265	10. PROGRAM ELEMENT, PROJECT, TASK 62702F 40270803
	12 REPORT DATE Aug 1979
Griffiss AFB NY 13441	13. NUMBER OF PAGES
MONITORING AGENCY NAME & ADDRESS(II different from Controlling Offi	ce) 15. SECURITY CLASS. (of this report) UNCLASSIFIED
Same 12/24P.	15a. DECLASSIFICATION/DOWNGRADING N/A
Approved for public release; distribution unlim	ited.
Approved for public release; distribution unliminated and the statement (of the abstract entered in Block 20, if different Same	
Approved for public release; distribution unlimitation un	nt from Report)
Approved for public release; distribution unlimited. 17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Richard S. Vonusa (IRAA) 19. KEY WORDS (Continue on reverse side if necessary and identify by block now Speech Processing	nt from Report)
Approved for public release; distribution unlimited. 17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Richard S. Vonusa (IRAA) 9. KEY WORDS (Continue on reverse side if necessary and identify by block now Speech Processing Digit Recognition Word Recognition Pattern Recognition	nt from Report)
Approved for public release; distribution unlimited. 17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different Same 18. SUPPLEMENTARY NOTES RADC Project Engineer: Richard S. Vonusa (IRAA) 19. KEY WORDS (Continue on reverse side if necessary and identify by block now Speech Processing Digit Recognition Word Recognition	mi from Report)) mber) lop the capability of recogtinuous speech, independent of ing this study contract was connected-digit recognition

DD 1 JAN 73 1473

UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE (When Data Entered)

347 650

m

SECURITY CLASSIFICATION OF THIS PAGE(When Data Entered)

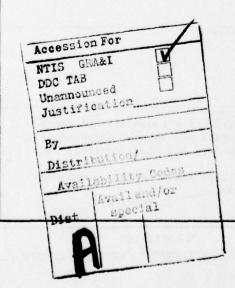
- (2) Clustering algorithms for use in the development of sets of reference patterns for speaker-independent word recognition.
- (3) Automatic enrollment for speaker-dependent, connected-word recognition for syntactically unconstrained word sequences of 20 words.

The program culminated in the installation of the speaker-independent, connected-digit recognition program on the Base and Installation Security System's advanced development model speaker verification system at RADG.

The speaker-independent, connected-digit recognition portion of this study resulted in a significantly faster algorithm with a 50-percent decrease in error rate over the course of this study-from 9.5 percent to 4.7 percent on an evaluation data set of 10 six-digit sequences from 106 speakers (64 males, 42 females).

The development of the clustering algorithm resulted in a two-stage, four-path algorithm with the mechanisms for detecting outlying data points in the design data and with subsequent analysis routines for comparing the results from the various paths and testing the validity of resulting clusters on the basis of comparisons with a priori information about the design data set.

The research into development of an automatic enrollment technique for speaker-dependent word recognition resulted in a method that yielded very good results for isolated word recognition but less acceptable results when used in continuous speech from the same speaker. The better results achieved with comparable hand enrollments point to the need for continued development of automated enrollment, and the desirability of an interim solution of a semiautomated enrollment procedure allowing the operator the option of modifying reference-point locations and recognition-pattern format definitions defined by an automated front end. Independent of the enrollment method, however, the benefit of reference file updating as a means of accommodating contextual variability, as well as intersession variability, was abundantly clear.



UNCLASSIFIED

SECURITY CLASSIFICATION OF THIS PAGE(When Date Entered)

EVALUATION

The objective of this program was to develop techniques and algorithms to extend highly reliable speaker-dependent isolated word recognition to speaker-dependent continuous word recognition and study the methodology for speaker-independent continuous speech recognition.

A hardware/software implementation of a real-time continuous speech recognition system was fabricated by Texas Instruments (TI). This system was extensively tested and modified to incorporate the results of the tests and a continual upgrade of the system took place over the life of the contract. TI based their real-time speech recognition system on techniques they developed for automatic speaker verification and the Total Voice Verification program which used a restricted connected digit capability.

The speaker-independent, connected digit performance resulted in 95.3 percent recognition accuracy on a data set consisting of 10 six-digit sequences from 106 speakers (64 males and 42 females).

The capability that TI has developed under subject program has been installed at RADC for further test and evaluation. The RADC tests shall attempt to establish the effectiveness of the current state-of-the-art connected speech system as to its applicability to operational military requirements.

Richard Conusa RICHARD S. VONUSA Project Engineer

TABLE OF CONTENTS

Secti	on Title	Page
1	INTRODUCTION	. 1
11	CONTINUOUS SPEECH RECOGNITION	. 3
	A. Speech Representation	. 10 . 15 . 17 . 19
Ш	SPEAKER-DEPENDENT WORD RECOGNITION	. 27
IV	A. Introduction	. 27 . 29
IV	WORD RECOGNITION	. 33
V	A. Review of Clustering in Speaker-Independent Word Recognition B. Detailed Clustering Algorithm C. Criteria for Measuring Partition Goodness D. Description of Cluster Analysis Documentation 1. Trees 2. Parameter Comparisons 3. Consistency Tests E. Observations on Clustering Results F. Testing Cluster Validity With a priori Information GENERAL-PURPOSE SPEECH I/O CAPABILITY	. 33 . 37 . 40 . 42 . 42 . 42 . 50 . 53
	A. System Description	. 76
VI	EXPERIMENTAL RESULTS	. 77
	A. Speaker-Independent Digit Recognition 1. Data Sets 2. Digit Recognition Results for Six-Digit Sequences 3. Digit Recognition Results for Three-Digit Sequences 4. Digit Recognition Results for Isolated Digits 5. Effect of Spectral Normalization Technique on Digit Recognition Performance 1. Digit Recognition Results for Isolated Digits 5. Effect of Spectral Normalization Technique on Digit Recognition	. 77 . 77 . 83 . 88
VII	B. Limited Vocabulary Word-Recognition Experiment	
	KEFEKERCES	

APPENDIXES

- Speech Processing
- Scatter Matrices B.
- C.
- tr S_B Alternate Form Derivation
 Postiterative Optimization Statistics for
 Recognition Patterns D.

LIST OF ILLUSTRATIONS

Figure	e Title	1	Page
1	Hierarchical Organization of Feature Abstraction Network H		3
2	Quantized Spectral Speech Representation		10
3	Demonstration of Need for Time Alignment Between Spectra for		
	Two "Sevens"		11
4	Relationship of Piecewise Linear Time Normalization to Nonlinear Time		
	Normalization		
5	Example of Recognition Pattern Formation		
6	Example of Scanning Pattern Formation		
7	Example of Valley Finding		
8	Example of Directed Graph		
9	Flow Chart for Efficient Tree-Searching Algorithm		
10	Word Recognition Algorithm Flow Chart		25
11	Automatic Enrollment for "Two"		30
12	Automatic Enrollment for "Two"		32
13	Bell Laboratories Clustering Procedures		35
14	Block Diagram of Clustering Procedures Developed for Clustering Scanning and		
	Recognition Patterns for Speaker-Independent Digit Recognition		38
15	Tree for Final 24 Stages of MINAVE Agglomerative Clustering of All (166)		
	Scanning Patterns for First Reference Point of Digit "Zero"		44
16	Parameter Comparisons for Pre- and Postiterative Optimization of MINAVE		
	Partitions Using All Points		45
17	Parameter Comparisons for Pre- and Postiterative Optimization of MINAVE		
	Partitions With Outliers Discarded		46
18	Parameter Comparisons for Pre- and Postiterative Optimization of MINMAX		
	Partitions Using All Points		47
19	Parameter Comparisons for Pre- and Postiterative Optimization of MINMAX		
	Partitions With Outliers Discarded		48
20	Parameter Comparisons for Preiterative Optimization of MINAVE and		
	MINMAX Partitions Using All Points		49
21	Parameter Comparisons for Postiterative Optimization of MINAVE and		
	MINMAX Partitions Using All Points		50
22	Class Assignments and Contingency Tables for Postiterative Optimization of		
	10 Classes Using All Points for Reference Point 1 of Digit 0		51
23	Contingency Table Measures for Pre- and Postiterative Optimization Using		
	All Points for Reference Point 1 of Digit 0		52
24	Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes		-
	for MINAVE Agglomerative Clustering of Scanning Patterns		54

25	Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes for MINAVE Agglomerative Clustering of Recognition Patterns			5
26	Histograms of Normalized Entropy as Measure of Dispersion in Class Size for Preiterative Optimization of Recognition Patterns			50
27	Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes for Preiterative Optimization of Recognition Patterns			
28	$\sigma_{\rm post}/\sigma_{\rm pre}$ Versus $\sigma_{\rm post}/\sigma_{\rm pre}$ for MINAVE Agglomerative Clustering of Scanning Patterns			
29	$\sigma_{ m post}/\sigma_{ m pre}$ Versus $\sigma_{ m post}/\sigma_{ m pre}$ for MINMAX Agglomerative Clustering of Scanning Patterns			
30	σ Versus α for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns			
31	σ Versus α for Preiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns			
32	σ Versus α for Postiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns			
33	Comparison of Mutual Information for Clustered Scanning Patterns Using Four Clustering Algorithms			
34	Comparison of Mutual Information for Clustering Recognition Patterns Using Four Clustering Algorithms			
35	Speech Channel Block Diagram			
36	Scanning Pattern Format			
37	Spectra for Digit "Two" for Speaker J.S			
	LIST OF TABLES			
Table	Title			
1	Title Characteristics of 16-Channel Filter Bank			9
1 2	Title Characteristics of 16-Channel Filter Bank			9
1 2 3	Title Characteristics of 16-Channel Filter Bank			9
1 2 3 4	Title Characteristics of 16-Channel Filter Bank			 91520
1 2 3 4	Characteristics of 16-Channel Filter Bank	 	 	 9 15 20 21
1 2 3 4 5	Characteristics of 16-Channel Filter Bank	 	 	 9 15 20 21 41
1 2 3 4 5 6	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization	 	 	 9 15 20 21 41 43 57
1 2 3 4 5 6 7 8	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization	 	 	 9 15 20 21 41 43 57
1 2 3 4 5 6 7 8 9	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns	 	 	9 15 20 21 41 43 57 58
1 2 3 4 5 6 7 8 9	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Postiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns	 	 	9 15 20 21 41 43 57 58
1 2 3 4 5 6 7 8 9	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns			9 15 20 21 41 43 57 58 66
1 2 3 4 5 6 7 8 9	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Postiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns			9 15 20 21 41 43 57 58 66 67 68
1 2 3 4 5 6 7 8 9	Characteristics of 16-Channel Filter Bank Recognition Pattern Format Definitions for Digits Sample of Information Contained in Table of Hypothesized Digits Sample of Table of Backpointers for Best Word Sequences in Sorted Table of Hypothesized Words Values of Seven Clustering Criteria for Postiterative Optimization of MINMAX Agglomerative Clusters for Recognition Pattern for Digit "Six" Criterion Values for Tree for Final 24 Stages of MINAVE Agglomerative Clustering for First Reference Point of "Zero" Contingency Table of First-Kind Results for Preiterative Optimization Contingency Table of First-Kind Results of Postiterative Optimization Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns Mutual Information and Residues for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns			9 15 20 21 41 43 57 58 66 67 68

Agglomerative Clusters of Recognition Patterns Mutual Information and Residues for Preiterative Optimization of MINMAX Agglomerative Clusters of Recognition Patterns Mutual Information and Residues for Postiterative Optimization of MINMAX Agglomerative Clusters of Recognition Patterns Ten Sets of Ten Six-Digit Sequences Used in Testing Martin-Herscher Digit Texts	 		. 7	72
Agglomerative Clusters of Recognition Patterns Mutual Information and Residues for Postiterative Optimization of MINMAX Agglomerative Clusters of Recognition Patterns Ten Sets of Ten Six-Digit Sequences Used in Testing	 		. 7	
Mutual Information and Residues for Postiterative Optimization of MINMAX Agglomerative Clusters of Recognition Patterns Ten Sets of Ten Six-Digit Sequences Used in Testing	 		. 7	
Agglomerative Clusters of Recognition Patterns				12
17 Ten Sets of Ten Six-Digit Sequences Used in Testing				72
and the second s				12
18 Martin-Herscher Digit Texts			. 7	78
			. 1	79
19 Synopsis of Evaluation Results			. 8	30
20 Digit Recognition Results for Selected Evaluation Runs			. 8	32
21 T-Function Quantization Thresholds			. 8	35
22 Decrease in Average Recognition Pattern Error by Including				
T-Function in Scanning Patterns			. 8	35
23 Confusion Matrix for Digit Recognition for 6-Digit Sequences for				
Run No. 62			. 8	35
24 Recognition Error (TE) Normalizing Constants			. 8	36
25 Confusion Matrix for Digit Recognition for 6-Digit Sequences for				
Run No. 73				
26 Digit-Recognition Performance of Males and Females				
27 Histogram of Digit-Recognition Performance			. 8	37
28 Confusion Matrix for Digit Recognition for 3-Digit Sequences				
Constrained in Length				
29 Confusion Matrix for Digit Recognition of Isolated Digits			. 8	39
30 Total Normalized Error (NE) for Digits From Sequence 852-734 for				
Speaker J.S.			. 9	1
Valley Point Errors, Sequence Errors (SQ), and Recognition Errors (TE)				
for Digits From Sequence 852-734 for Speaker J.S				
32 σ_{\min} Variation Performance Test				
33 Confusion Matrix for Automatic Enrollment				
34 Confusion Matrix for Automatic Enrollment			. 9)5
35 Percent Correct Speaker-Dependent Recognition Results for Continuous				
Utterances From a 20-Word Vocabulary Enrolled on Isolated Words			(16

SECTION I

This final report covers the research done on a limited-vocabulary continuous word recognition study undertaken by Texas Instruments. This effort was divided into two primary areas of investigation: extension of speaker-dependent isolated word recognition to speaker-dependent continuous word recognition, and the study of speaker-independent continuous speech recognition.

Speaker-dependent isolated word recognition is currently being used for applications such as map data entry. Extension to speaker-dependent continuous word recognition is a more natural one for the time normalization techniques used at Texas Instruments (described in Section II) than the type that depends on locating endpoints of words, which may not even exist [e.g., when phonemes are shared (/s/ in six-seven) between words in continuous speech]. Speaker-dependent word recognition uses speaker-dependent reference patterns obtained in a single enrollment session. A method of automatic enrollment and supervised updating to accommodate intersession variations and context dependencies were investigated during this study.

For many years, the approach to the problem of speaker-independent recognition of continuous speech has been a heuristically directed search for the correct features and weightings for the hierarchical classification of a set of symbol strings, mapping ultimately into an English-language transcription. The emphasis has been on getting out of the acoustic and into the phonemic domain as quickly as possible because of the huge memory requirements for storing acoustic data for large vocabularies. Since the heuristics were often based on the researcher's judgment, derived from often insufficient data, the consequent mislabelings had to be corrected with progressively more complex classification algorithms. Design and testing using small data bases, along with the use of phonemic representations of speech have resulted not only from memory limitations but also from the lack of techniques in speech for dealing with very large amounts of data. Within the last few years, however, work on such techniques has begun to appear. During the Total Voice Speaker Verification study, performed by Texas Instruments under RADC sponsorship, a clustering algorithm was developed and used to produce a set of speaker-independent reference patterns for use in speaker-independent, connected-digit recognition. The current study then concentrated on two tasks for speaker-independent, continuous-speech recognition. One task was to determine the performance that could be achieved on speaker-independent, connected-digit recognition using the previously developed reference patterns by making improvements to the word recognition algorithm. The other task was to investigate improvements that could be made to the clustering algorithm for the purpose of finding better partitions of the design data set.

Section II of this report reviews the speech technology used at Texas Instruments and covers an improved directed graph searching algorithm developed during this contract. Section III covers an automated enrollment method for speaker-dependent, connected-word recognition and the role of reference-pattern updating. Section IV describes the application of clustering to speaker-independent reference-pattern generation and covers the algorithm extensions developed

¹ R.L. Davis, B.M. Hydrick, and G.R. Doddington, "Total Voice Speaker Verification," Rome Air Development Center Technical Report, RADC-TR-78-260, January 1978.

during this contract. Section V covers several general-purpose speech-processing capabilities that center on the use of direct speech input/output (I/O) to a fast array processor. The experimental results for both the extensive testing performed for speaker-independent, connected-digit recognition and the more limited testing done for speaker-dependent, continuous-speech recognition are covered in Section VI. Conclusions and recommendations are made in Section VII.

SECTION II CONTINUOUS SPEECH RECOGNITION

During the relatively short history of continuous speech recognition work, the classification schemes have used a feature abstraction process from the speech waveform followed by a hierarchical classifier. The level of abstraction varied from features of the waveform itself to symbol representations (phonemes) requiring highly sophisticated classification techniques in order to compensate for segmentation and labeling errors. The classification complexity generally was proportional to the level of abstraction. Martin² shows a tree (Figure 1) of feature abstraction levels.

The usual argument for using a symbol is that it offers a more compact representation of words and, hence, growth in the memory requirement is not so dramatic with increase in vocabulary size. However, as Reddy³ points out, good signal-to-symbol transformation techniques currently do not exist, causing size increases in the lexicons and the algorithms, not only to account for context, dialect, and idolect variations, but also to account for mislabeled acoustic events.

Therefore, reference-pattern matching in the signal domain has the advantage of not having to accommodate feature abstraction errors. Three crucial problems are involved, however: selection of the speech representation, time normalization of the speech signal for matching with reference patterns, and selection of the reference patterns themselves. The first two of these topics are discussed in the remainder of this section and the reference-pattern selection is the subject of the more extended discussions in Sections III and IV.

A. SPEECH REPRESENTATION

The specific speech representation used in this study was the output of a 16-channel digital filter bank preceded by a first-order differencing network (for preemphasis). Each of the bandpass filters is a two-section, cascaded, second-order Bessel filter followed by a rectifier and a lowpass filter sampled every 10 milliseconds (ms). Center frequencies and bandwidths for the 16 filters are given in Table 1.

An important consideration is the choice of wide-bandwidth filters that locate spectral peaks but that avoid resolving the voice fundamental and its harmonics. Note that the center frequencies for filters 14 through 16 lie in the part of the spectrum primarily occupied by energy only during fricatives. The exception is the third formant for the vowel /i/ for males and for all the front vowels for females. Since no precise resolution of the frequency location is possible with the wide-bandwidth filters, the only interest is the presence or absence of a third formant (in which case other formants would also exist in lower filters) or the presence or absence of energy anywhere in the frequency band of the top three filters without lower frequency energies. In order to compact the filter bank representation, the top three filters were added into one value, without averaging because of the depressed amplitudes, yielding a 14-element vector to represent the speech spectrum:

April 1976.

²T.B. Martin, "Acoustic Recognition of a Limited Vocabulary in Continuous Speech," Ph.D. Dissertation, University of Pennsylvania, 1970.

³D.R. Reddy, "Speech Recognition by Machine: A Review," *Proceedings of the IEEE*, 64:501-531,

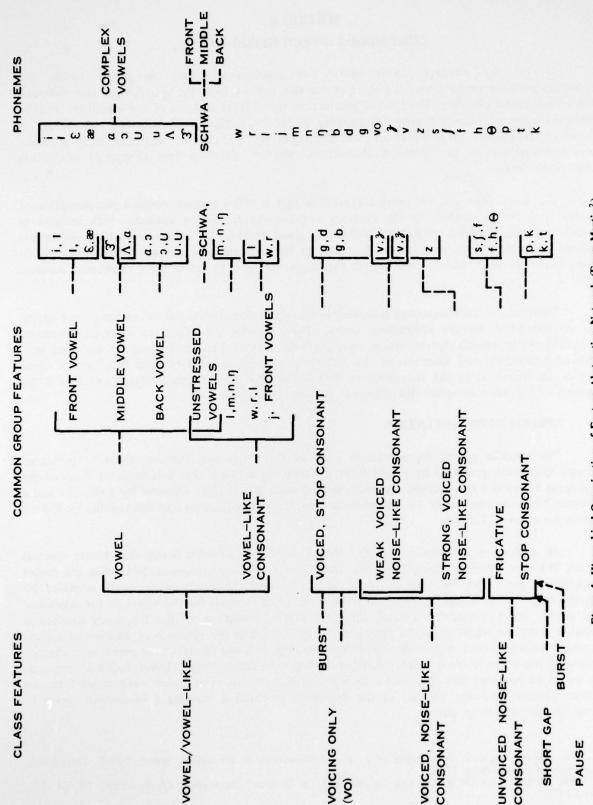


Figure 1. Hierarchical Organization of Feature Abstraction Network (From Martin²)

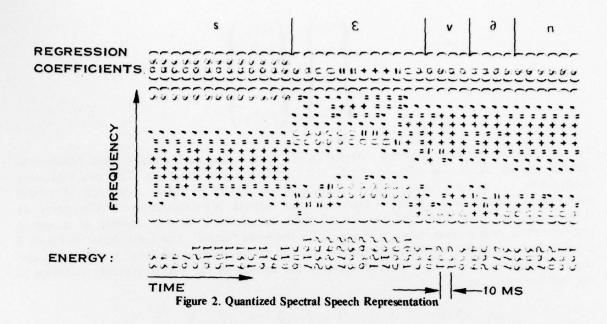
$$\overrightarrow{A}_{j} = \left\{ \begin{array}{c} a_{1j} \\ a_{2j} \\ \vdots \\ a_{14j} \end{array} \right\} = \left\{ \begin{array}{c} a_{1}(t_{j}) \\ a_{2}(t_{j}) \\ \vdots \\ a_{14}(t_{j}) \end{array} \right\}$$

This 14-element vector is regressed (Appendix A) using a sine and a cosine basis function to eliminate gross aspects of the spectrum and to flatten the spectrum. The two regression coefficients, c_1 and c_2 , along with a measure of the energy in filters 2 through 13 (vowel energy), are concatenated to the 14 regressed filter outputs. All elements except the energy are then normalized and quantized to one of eight equiprobable values, resulting in a speech representation such as that shown in Figure 2 for the word "seven." The form shown in Figure 2 is used throughout the remainder of this report. The values of the normalized, quantized a_{ij} and the two regression coefficients are indicated by the density of the printed symbols according to the following:

At this point the energy is not quantized; however, it is always used relative to other energies and the relative value is then quantized. Further detail of the speech representation can be found in the Total Voice Speaker Verification study final report.¹

TABLE 1. CHARACTERISTICS OF 16-CHANNEL FILTER BANK

Filter	Center Frequency (Hz)	Bandwidth (Hz, at -6 dB)
1	280	250
2	395	280
3	525	310
4	630	340
5	750	360
6	900	360
7	1080	360
8	1265	365
9	1480	365
10	1725	365
11	1985	365
12	2285	360
13	2640	365
14	3150	625
15	3720	635
16	4235	615



B. TIME NORMALIZATION

One of the basic problems in speech processing is time alignment of the speech waveform with respect to a reference. For example, in the two spectrograms in Figure 3 for the word "seven," the time differences between corresponding Δs (which denote phonemic boundaries) are obvious.

Early work used linear time normalization of two patterns between endpoints of words, and although this method improved recognition performance, it suffered from an inability to deal with the nonlinear fluctuations between endpoints and to locate endpoints in continuous speech reliably.

Two distinct approaches developed during the late 1960s and early 1970s. One approach (most of the ARPA sponsored work: Reddy³) was based on translating a string of input features into a sequence of phonemic labels, a procedure dependent on accurate segmentation between phonemes. Segmentation and labeling errors were then repaired by more sophisticated subsequent processing using syntax, semantics, etc.

The other method approached the problem by a nonlinear warping of the time axis of a feature waveform of the input speech to obtain maximum coincidence with a reference waveform. This approach was used by both Doddington⁴ and Sakoe and Chiba;⁵ however, the latter's approach could be more easily represented in a form amenable to the use of dynamic programming, useful in easing the computation burden. This dynamic programming approach has

⁴G.R. Doddington, "A Method of Speaker Verification," Ph.D. Dissertation (Thesis), University of Wisconsin. 1970.

⁵H. Sakoe and S. Chiba, "A Dynamic Programming Approach to Continuous Speech Recognition," Proceedings of the 7th International Congress on Acoustics, August 1971.

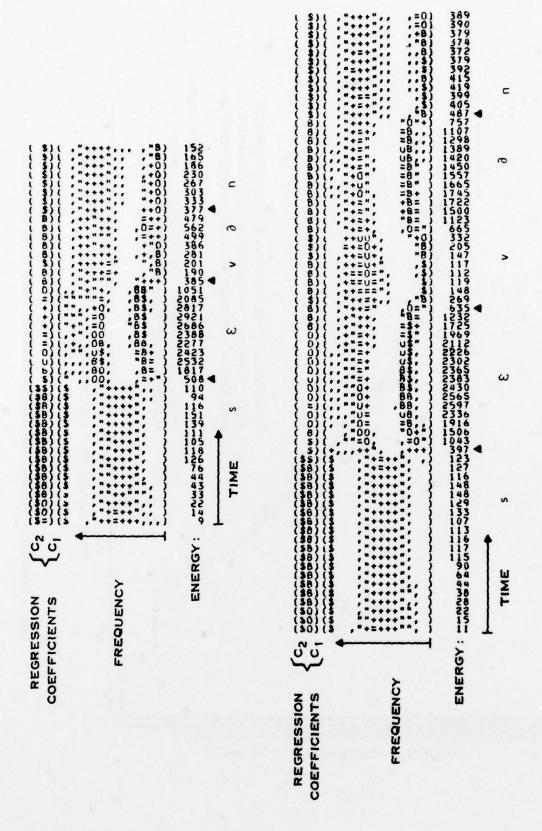


Figure 3. Demonstration of Need for Time Alignment Between Spectra for Two "Sevens"

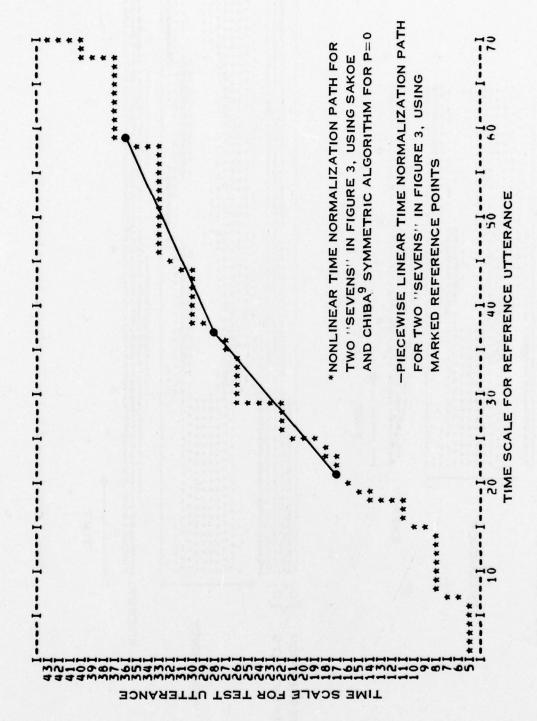


Figure 4. Relationship of Piecewise Linear Time Normalization to Nonlinear Time Normalization

been used by Velichko and Zagoruiko,6 Itakura,7 White and Neely,8 and Sakoe and Chiba9 in isolated word recognition. Extension to continuous-word recognition has been done by Lowerre¹⁰ on the HARPY speech recognition system (also described briefly by White¹¹), by Porter. ¹² and by Nippon Electric Company in their DP-100 Connected Speech Recognition System.

The technique used at Texas Instruments is an amalgamation of these two methods and was first used by Doddington¹³ in 1973. In this method, potential acoustic boundaries (reference points) are first located in the input waveform. Reference points are combined into optimal sequences for words in the vocabulary using a dynamic programming routine that uses a measure of how reliably the reference points were located and that accounts for the deviations from expected time differences between reference points.

After sequences of potential reference points have been identified, the input waveforms are interpolated linearly between reference points to form a time normalized representation of the utterance. The relationship to the Sakoe/Chiba approach can be seen in Figure 4. Essentially, only those points along the path of the time warp that represent acoustic boundaries are found (the o's in Figure 4), and the linear interpolation is then performed between these reference points. A piecewise linear, time-normalized, acoustic representation of the word (a "recognition pattern") is thus formed. A sample of the spectral data portion of a recognition pattern being extracted from input speech spectra is shown in Figure 5.

As an example of the choice of reference-point locations, the reference points (Δs) for the digits are shown in Table 2 for the phonetic transcriptions of the general American dialect pronunciations for the digits as found in Kenyon and Knott.¹⁴ These locations were chosen at points that would exhibit large spectral changes. The actual rules used in extracting recognition patterns for the 10 digits are also specified in table 2, where:

- Initial negative numbers indicate the columns for extrapolation before the first (1) reference point
- Intermediate numbers in parentheses indicate the number of columns for (2) interpolation between reference points
- The remaining numbers indicate columns for extrapolation after the last reference (3)
- ⁶V.M. Velichko and N.G. Zagoruiko, "Automatic Recognition of 200 Words," International Journal Man-Machine Studies, 2:223, June 1970.
- ⁷F. Itakura, "Minimum Prediction Residual Principle Applied to Speech Recognition," IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-23:67-72, February 1975.
- ⁸G.M. White and R.B. Neely, "Speech Recognition Experiments With Linear Prediction, Bandpass Filtering, and Dynamic Programming," IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-24:183-188,
- April 1976.

 "H. Sakoe and S. Chiba, "Dynamic Programming Algorithm Optimization for Spoken Word Recognition,"
- IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-26:43-49, February 1978.

 10 B.T. Lowerre, "The HARPY Speech Recognition System," Ph.D. Dissertation (Thesis), Carnegie-Mellon
- University, 1976.

 11 G.M. White, "Continuous Speech Recognition: Dynamic Programming, Knowledge Nets and HARPY," Paper 28-2, 1978 WESCON Professional Program, September 1978.
- ²J.E. Porter, "LISTEN: A System for Recognizing Connected Speech Over Small, Fixed Vocabularies in Real
- Time," Naval Training Equipment Center Technical Report, NAVTRAEQUIPCEN 77-C-0096-1, April 1978.

 13 G.R. Doddington, "Speaker Verification," Rome Air Development Center Technical Report, RADC-TR-74-179, April 1974.
- ¹⁴ J.S. Kenyon and T.A. Knott, A Pronouncing Dictionary of American English, G. & C. Merriam Company (Springfield, Massachusetts, 1953).

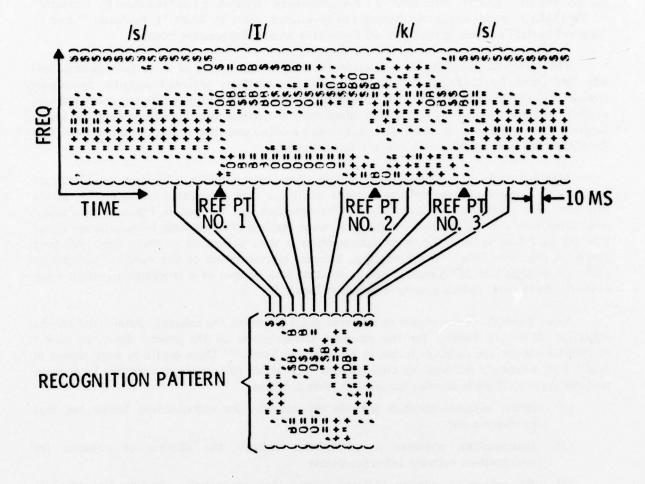


Figure 5. Example of Recognition Pattern Formation

At this point, the speech representation is still in the acoustic domain, differing from those 10,15 who transform their time-warped segments into phonemic labels with associated transition probabilities between labeled states. The advantage of remaining in the acoustic domain is that it avoids an intermediate classification that would introduce errors and obviates the need to find every phonetic boundary, which is helpful when such boundaries are difficult to find.

¹⁵F. Jelinek, "Continuous Speech Recognition by Statistical Methods," *Proceedings of the IEEE*, 64:532–556, April 1976.

TABLE 2. RECOGNITION PATTERN FORMAT DEFINITIONS FOR THE DIGITS

C. REFERENCE-POINT LOCATION

A presupposition of this piecewise, linear, time-normalization technique is extremely accurate reference-point location. One approach would be vocabulary-independent, locating changes in features such as voicing, energy, or spectrum between adjacent time samples. This is a reliable, precise method for use in speaker-dependent recognizers; however, sometimes expected acoustic segmentation points are missed in speaker-independent recognition.

A more robust approach is to use a vocabulary-dependent approach (similar to the "transeme" approach used at IBM¹⁶), matching a feature vector (called a "scanning pattern") extracted from the input speech waveform to reference scanning patterns, or templates. Figure 6 shows a scanning pattern being extracted from the spectral input. Matching is performed by computing a distance between the input and all reference patterns for every frame (10 ms in this study). Minima in this distance function are locations of potential acoustic boundaries (reference points).

More specifically, the scanning pattern formed at time t_j consists of: (1) the spectral data, regression coefficients, and energy for the five time samples from t_{j-2} through t_{j+2} and (2) the difference between the data for all adjacent pairs of time samples. The energy used in the scanning pattern is the energy for each of the five columns of data, normalized by the sum over all five columns and quantized to 4 bits. Figure 6 illustrates the formation of a scanning pattern from preprocessed (regressed, normalized, quantized) speech data. The only purpose of the difference data is to weight more heavily rapid changes of the feature vectors with respect to time. Since these data are derived from the standard data portion of the scanning pattern, subsequent illustrations of scanning patterns in this report will not show the difference data, even though it is, in fact, part of the actual pattern.

¹⁶N.R. Dixon and H.F. Silverman, "The 1976 Modular Acoustic Processor (MAP)," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-25:367-379, October 1977.

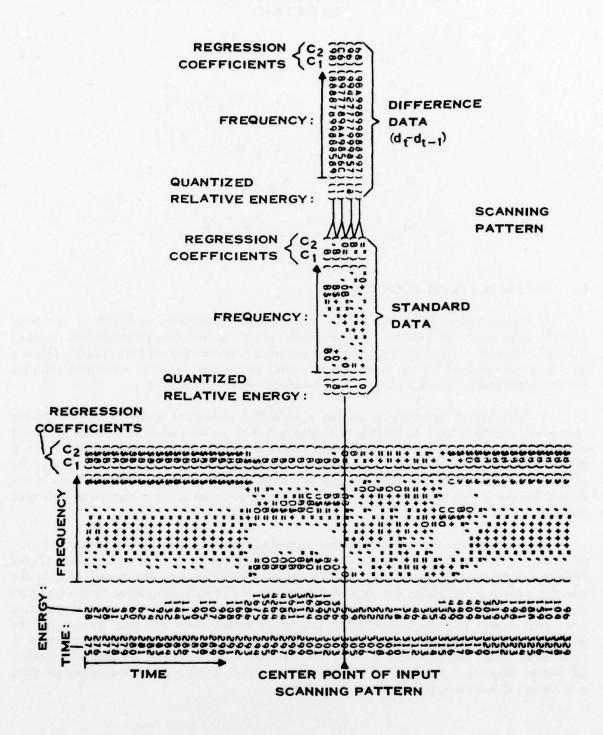


Figure 6. Example of Scanning Pattern Formation

In order to determine where reference points occur in the input speech, the input data are compared with reference data. This procedure (called scanning) is done by formatting scanning patterns from the input speech at each time sample t_j , comparing these with predetermined reference scanning patterns \overline{t}_k , and obtaining a measure of squared difference between the two, called the scanning error:

$$e_{kj} = ||\vec{x}_j - \vec{r}_k||^2 = \sum_{i=1}^{164} (x_{ij} - r_{ik})^2$$
 (1)

The final error associated with each reference point is the minimum error of all comparisons with patterns representing that reference point.

Using the scanning errors as a function of time, an error function is thus generated for each type of reference scanning pattern using the minimum scanning error for each pattern type for each time sample. (Multiple reference scanning patterns may be allowed for each reference point of each word.) Each function is monitored for dips of sufficient magnitude to be considered as potential locations of the corresponding reference points in the input data. These dips are called valley points when the ratio of the scanning error following the dip to the scanning error at the dip itself is greater than or equal to a specified peak-to-valley ratio (PVR), typically 1.1 to 1.3, and the magnitude of the scanning error for the valley point is less than or equal to a threshold, typically 600 to 1,200. The occurrence of a peak (verified when the ratio of the scanning error following the peak to the scanning error at the peak is less than the reciprocal of the PVR) is required before another valley point can be found. The valley-finding procedure is shown in Figure 7.

D. WORD HYPOTHESIZING AND TESTING

Once these valleys in the scanning error (potential reference points) have been found, the next task is to fit them together to form word hypotheses. A sequence of time-ordered reference-point hypotheses for a word must exist, and the time distance between each pair of reference points must satisfy word-specific minimum/maximum restrictions. The error determined for each reference point pair is weighted by deviations from the expected distance between the two points and the scanning error at each hypothesized reference point. The weighted error for reference points i and j is:

$$E_{w_{i,j}} = \frac{(e_i + offset) (e_j + offset)}{1024} \left[1 + \beta \left(\frac{\hat{d}t_{i,j} - dt_{i,j}}{\hat{d}t_{i,j}^*} \right)^2 \right]$$

where

$$dt_{i,j} = t_j - t_i \qquad \beta = 2$$

$$dt_{i,j} = \text{expected } dt_{i,j} \qquad dt_{\min} = 4$$

$$dt_{i,j}^* = \max (dt_{i,j} dt_{\min}) \qquad \text{offset } = 100$$

$$e_i, e_i = \text{scanning error for reference points } i, j$$

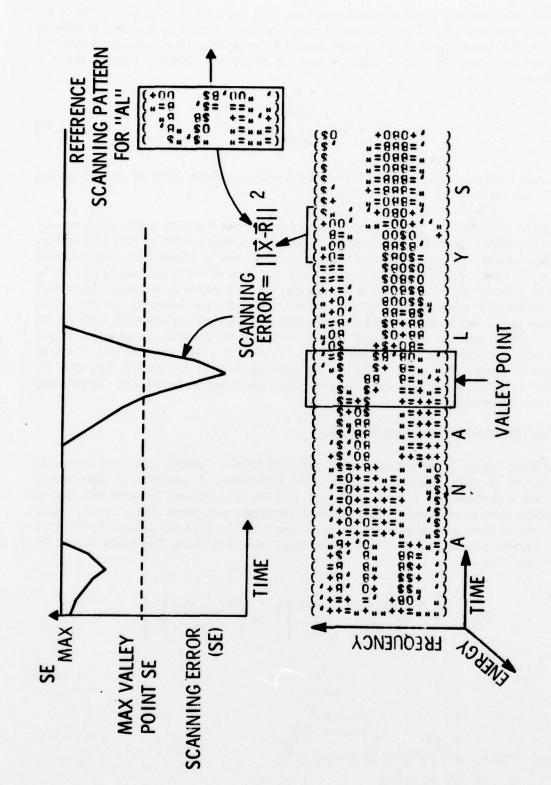


Figure 7. Example of Valley Finding

If the hypothesized word sequence error (SQ), which is the sum of the $E_{wi,j}$ for all reference point pairs in the word, is less than a predetermined word-specific threshold, then a word has been hypothesized.

To test this hypothesized word, the time-normalized recognition pattern anchored at the corresponding hypothesized reference points is compared to time-normalized reference recognition pattern(s) for the hypothesized word, using the squared Euclidian distance. This distance (or the minimum distance, in the case of multiple reference patterns) is the recognition error (TE), which is used along with the sequence error (SQ) in computing a total normalized error (NE) for the word "k" as given below:

$$NE_{k} = \frac{TE_{k}/No. \text{ of columns in word } k}{\text{normalizing constant for word } k} + w_{k} \frac{SQ_{k}}{10 \text{ (NPP)}}$$
(3)

where NPP is the number of reference-point pairs used in computing SQ_k . If the NE for the hypothesized word is above a predefined threshold or if the average energy across the recognition pattern is less than a threshold, the word is discarded. Otherwise, it is placed into a table of hypothesized words, along with the SQ, TE, NE, average energy, reference-point times, and scanning errors for that word. Table 3 is an example table of hypothesized words for a sequence of six digits spoken continuously. Note the existence of the 35 superfluous words in this case. The optimal searching of this table (or directed graph) is described in the next subsection.

E. AN EFFICIENT TREE-SEARCHING ALGORITHM

Once the current speech segment has been completed, the sorted table of hypothesized words must be searched to find the sequence having the minimum error. Note that this corresponds to finding the best path through a directed graph, such as the simple one shown in Figure 8, taken from Porter. 12 The correspondence to the directed graph can be seen most readily from a specific example. Table 4 shows the table of hypothesized words (digits, in this case) sorted according to the time (in centiseconds) of the final reference points. Each entry in this table can point back in time to all previous digits with earlier final reference-point times (less than the initial reference-point time of the entry being considered). The allowable range of backpointers can be limited by requiring that the time difference between the first reference point of each word and the last reference point of preceding words lie between a specified minimum and maximum.

Clearly, the exhaustive search time for traversing all possible paths increases rapidly.

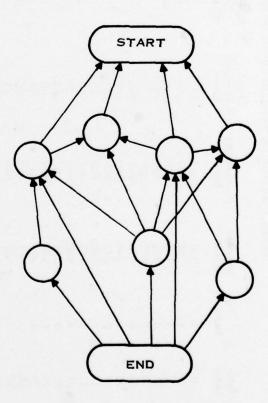


Figure 8. Example of Directed Graph

	Ref. Pt. 3 Time Error				191										382			386		519		139			220		231			112				312	338	273					
	Ref. Time				469										412			423		413		392			400		392			362				329	353	357					
	Ref. Pt. 2 ime Error	448	180	623	296	493	493	477	423	423	394	277	277	319	254	332	265	405	260	187	219	263	229	215	451	240	230	366	358	324	318	318	377	313	313	172	205	589	124	348	187
	Ref. Time	454	469	454	455	435	435	436	430	430	430	422	422	423	391	423	423	408	422	394	400	381	399	398	373	394	386	351	358	355	359	359	355	330	330	343	313	268	271	271	271
DIGITS	t. 1 Error	295	345	382	501	272	361	295	270	295	344	272	399	427	256	109	344	438	354	405	463	256	638	427	251	425	405	246	323	263	329	422	329	300	300	437	345	396	290	392	426
HESIZED	Ref. Pt. 1 Time Error	415	439	439	439	394	417	415	396	415	409	394	401	401	373	405	409	404	411	374	387	373	374	374	369	376	374	343	332	331	331	342	331	325	325	332	291	250	255	250	248
LE OF HYPOT	Average Energy	1,283	908	1,173	1.148	646	176	875	595	916	988	787	785	736	870	756	1,168	854	1,145	840	510	1.094	784	1,000	672	1,011	1,116	2,054	1,464	952	1,406	1,788	1,586	1,043	1,240	1,622	649	1,120	1,664	1,067	1,180
SAMPLE INFORMATION CONTAINED IN TABLE OF HYPOTHESIZED DIGITS	Normalized Error Threshold	114	113	110	128	114	114	114	114	114	110	114	114	114	109	113	110	123	97	128	107	109	113	110	123	6	128	114	107	110	110	. 011	110	123	123	128	113	114	107	113	110
	Total Normalized Error	112	080	66	77	106	108	103	96	66	73	88	06	86	901	113	28	100	82	95	64	53	111	20	75	96	99	96	82	83	65	8/2	7.2	19	70	83	51	91	53	66	83
	Recognition Error	1,055	315	511	868	951	993	866	940	931	405	206	888	932	857	994	337	625	811	734	361	487	1,002	408	464	656	482	823	503	531	381	462	395	381	446	029	446	888	353	916	537
TABLE 3	Sequence Error	522	447	716	495	576	544	446	400	452	450	274	368	438	993	602	346	857	318	631	374	324	532	394	571	362	542	578	390	494	360	428	486	535	563	472	270	384	170	432	340
	Word	6	0 4	t ("	0	6	6	6	6	6	3	6	6	6	7	2	3	7	-	0	00	7	'n	3	7	-	0	6	00	9,		· ·	n (2	2	0	2	6	∞	2	m
	able	4 5	9 0	38	37	36	35	34	33	32	31	30	56	28	27	56	25	24	23	22	71	20	61	18	17	16	15	4 3	2	17	= :	0	70	× 0	7	9	2	4	3	7	-

TABLE 4. SAMPLE OF TABLE OF BACKPOINTERS FOR BEST WORD SEQUENCES IN SORTED TABLE OF HYPOTHESIZED WORD

Sorted Table	Time of First	Time of Last	Adjusted Normalized			Back point Sequence			
Index	Ref. Pt.	Ref. Pt.	Error	1	2	3	4	5	6
Best Sequ	uence:			16	16	16	16	28	39
41	439	469	80	-1	16	16	16	16	28
40	439	469	69	-1	16	16	16	16	28
39	439	464	59	-1	16	16	16	16	28
38	415	454	112	-1	16	16	16	16	-1
37	439	454	99	-1	16	16	16	16	28
36	415	436	103	-1	16	16	16	16	-1
35	394	435	106	-1	5	5	11	-1	-1
34	417	435	108	-1	16	16	16	16	-1
33	396	430	96	-1	16	16	16	16	-1
32	415	430	99	-1	16	16	16	16	-1
31	409	430	73	-1	16	16	16	16	-1
30	401	423	98	-1	16	16	16	16	-1
29	405	423	113	-1	16	16	16	16	-1
28	409	423	58	-1	16	16	16	16	-1
27	404	423	90	-1	16	16	16	16	-1
26	394	422	88	-1	5	5	11	-1	-1
25	401	422	90	-1	16	16	16	16	-1
24	411	422	82	-1	16	16	16	16	-1
23	374	413	85	-1	5	5	11	-1	-1
22	373	412	95	-1	5	5	11	-1	-1
21	387	400	64	-1	5	5	11	-1	-1
20	369	400	67	-1	5	5	11	-1	-1
19	374	399	111	-1	5	5	11	-1	-1
18	374	398	70	-1	5	5	11	-1	-1
17	376	394	96	-1	5	5	11	-1	-1
16	373	392	47	-1	5	5	11	-1	-1
15	374	392	59	-1	5	5	11	-1	-1
14	331	362	74	-1	5	5	-1	-1	-1
13	331	359	64	-1	5	5	-1	-1	-1
12	342	359	78	-1	5	5	-1	1	-1
11	325	359	54	-1	5	5	-1	-1	-1
10	332	358	82	-1	5	5	-1	-1	-1
9	332	357	74	-1	5	5	-1	-1	-1
8	331	355	72	-1	5	5	-1	-1	-1
7	325	353	63	-1	5	5	-1	-1	-1
6	343	351	96	-1	5	5	-1	-1	-1
5	291	313	51	-1	4	-1	-1	-1	-1
4	255	271	53	-1	-1	-1	-1	-1	-1
3	250	271	99	-1	-1	-1	-1	-1	-1
2	248	271	83	-1	-1	-1	-1	-1	-1
1	250	268	91	-1	-1	-1	-1	-1	-1

Syntactic constraints, such as those used in the total voice speaker verification study, can sometimes aid the efficiency of sequence finding, depending on whether the constraints apply locally to word pairs or more globally to the utterance as a whole. If syntactic constraints have been imposed in order to increase sequence recognition performance, such as in the total voice application, these constraints can be incorporated into the tree-searching algorithm to eliminate searching branches that are not syntactically correct.

An even greater potential saving may be obtained by saving optimal subsequences to avoid repetitious searching of the same path. However, this technique is not acceptable for sequences that are syntactically constrained since the saved optimal subsequence may not, in [fact, satisfy the constraints when those constraints are applied to the entire sequence. In other words, the "correct" sequence (satisfying the syntactic constraints) may in fact not be the "best" (lowest error) sequence. Since the saving of optimal subsequences only finds one best sequence of a given length, this method is not appropriate to syntactically constrained word sequences.

For unconstrained sequences, however, this technique of saving optimal subsequences can shorten the exhaustive search time to be proportional only to the number of table entries preceding a table entry. As an example, Table 4 gives the resulting subsequence backpointers for the sorted table of hypothesized digits for the six-digit sequence in Table 3. These backpointer lists are constructed from the bottom up, according to the algorithm shown in Figure 9. Essentially, the backpointer for sorted table index i, sequence length k, points to the sorted table index j (j \leq i) that has the minimum error for a subsequence of length (k - 1) for all j satisfying the relation:

$$\Delta t_{\min} \le (t_i^{\text{initial}} - t_j^{\text{final}}) \le \Delta t_{\max}$$
 (4)

where $\Delta t_{min} = 3$ cs and $\Delta t_{max} = 120$ cs were used in this study.

As each new entry in the list of backpointers is constructed, it is compared to the best (lowest total error) subsequence of the same length. If the newer entry has a lower error, this error replaces the best error, and the pointer to the best sequence of that length is changed to point to the sorted table entry currently under consideration. After the final sorted table index has been completed, the array of pointers to the subsequences having the lowest error contains the optimal results of the search. If the length of the sequence has been constrained, all that is necessary is to select the backpointer for the specified length sequence. If not, the second half of the algorithm shown in Figure 9 is used to determine the best sequence out of those specified by the array of pointers to the best sequence.

Although the algorithm described in this subsection was developed as a natural extension to existing sequence finders that did not save optimal subsequences, reference should be made to the work of others on this problem. Most widely published is the work of Jelinek et al. at IBM. 15,17,18 The IBM work was predicated on a phonemic representation of the recognized speech. The descriptions were in terms of probabilistic finite state machines where the recognized phonemes are outputs of state-to-state transitions, with all state transitions having associated a priori probabilities.

¹⁷F. Jelinek, L.R. Bahl, and R.L. Mercer, "Design of a Linguistic Statistical Decoder for the Recognition of Continuous Speech," *IEEE Transactions on Information Theory*, IT-21:250–256, May 1975.

¹⁸L.R. Bahl and F. Jelinek, "Decoding for Channels With Insertions, Deletions, and Substitutions With Applications to Speech Recognition," *IEEE Transactions on Information Theory*, IT-21:404–411, July 1975.

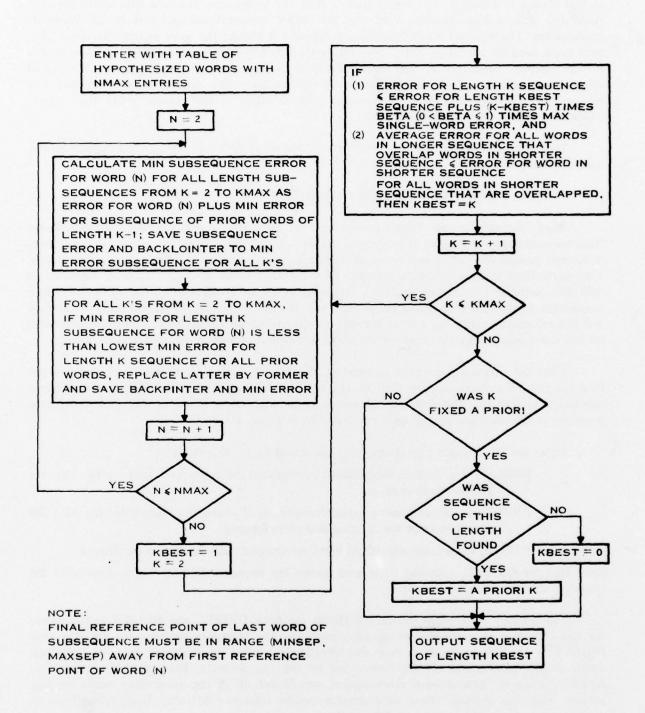


Figure 9. Flow Chart for Efficient Tree-Searching Algorithm

Much closer to the work of this subsection is that of Porter, 12 whose MINT algorithm is used to find the highest probability path through a directed graph of hypothesized words, such as that shown in Figure 8. The nodes (rather than the transitions, as in the IBM work) are each associated with a hypothesized word and the edges represent backpointers to all allowable predecessors. The solution to the real-time processing is almost the same as that described here; each node need be processed only once, rather than each sequence once, by saving backpointers to the optimal subsequence. The only difference is that, in the present study, N optimal subsequences are saved (N is the maximum allowable length sequence), whereas just one backpointer is saved by Porter. Subsequent postprocessing in the present study then selects which of the final N sequences is the best.

The interested reader is referred to Porter¹² for an extended discussion of the probabilistic basis for this procedure.

F. OVERALL WORD-RECOGNITION ALGORITHM

Word recognition at Texas Instruments is currently based on the piecewise-linear time-normalization technique (Subsection II.B) of finding potential acoustic boundaries (reference points) and fitting sequences of reference points together to form hypothesized words. Time-normalized spectral patterns, formatted for the hypothesized words, are then compared to reference patterns (for either speaker-independent or speaker-dependent recognition of either continuous or discrete speech). If the comparison is a good enough match, the time of the first and last reference points and a total normalized error (distance between the input and reference) for the word, along with the label for the word, are stored in a table of hypothesized words.

After the utterance has been completed, the table is sorted by time of occurrence of the final reference point and is then used by the tree-searching algorithm described in the previous subsection to find the best sequence of words. A summary flow chart for the word recognition programs investigated during this contract is shown in Figure 10.

Three specific computer programs were generated during this study:

DIGREC, for speaker-independent recognition of connected digits using only the TI 980 minicomputer

DIGRCT, for speaker-independent recognition of connected digits using the AP 120B array processor for filtering and preprocessing

RTENR, for speaker-dependent word recognition with automatic enrollment.

Note that for DIGREC, sampling is stopped during the sequence finding that is done after the complete utterance has been input.

The primary differences between RTENR and both DIGREC and DIGRCT is the source for the reference scanning and recognition patterns. The reference patterns for DIGREC and DIGRCT were derived from a clustering procedure applied to a design data set collected, digitized and hand-edited off-line before use by the test subjects. The reference patterns for RTENR, however, were derived from on-line enrollment of all the vocabulary words by each subject using the system. These were speaker-specific reference patterns. Both procedures are described in more detail in Sections III and IV.

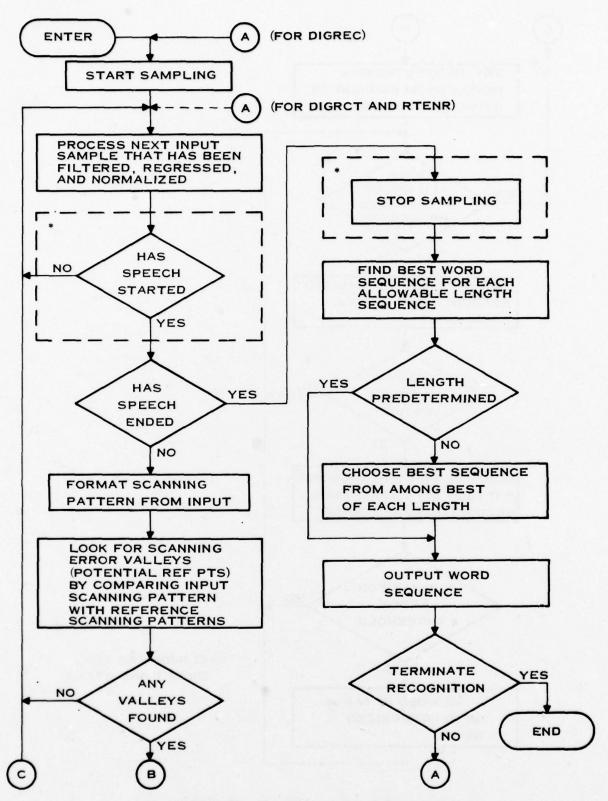


Figure 10. Word Recognition Algorithm Flow Chart (Sheet 1 of 2)

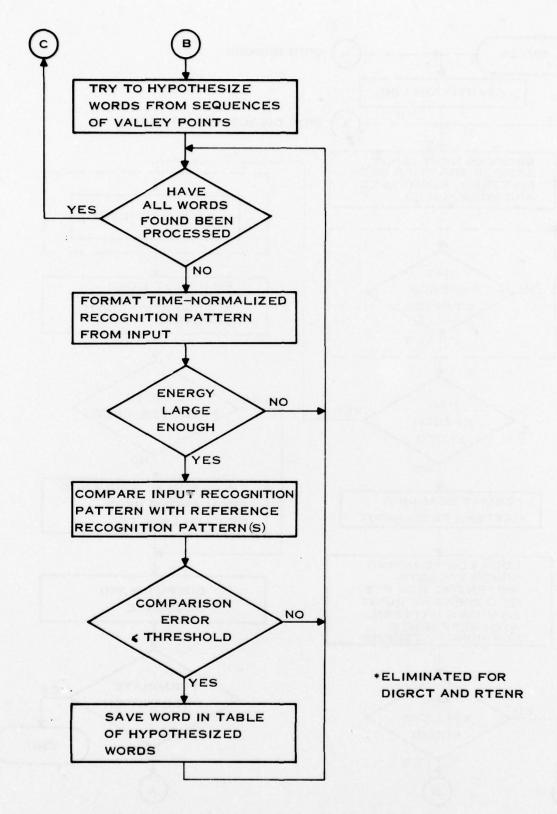


Figure 10. Word Recognition Algorithm Flow Chart (Sheet 2 of 2)

SECTION III SPEAKER-DEPENDENT WORD RECOGNITION

A. INTRODUCTION

The speech-processing algorithm described in Section II provides the framework for both speaker-dependent and speaker-independent word-recognition tasks. The algorithm is made specific to the given task through definition of the reference scanning and recognition patterns. The speaker-dependent task is accomplished through definition of a single set of reference scanning patterns for each vocabulary word for each speaker. In contrast, multiple scanning patterns are defined for each word in the speaker-independent word-recognition task (discussed in Section IV).

The reference patterns for the speaker-dependent word-recognition task are obtained in a single enrollment session where each word in the vocabulary is spoken in isolation. Intersession variations and contextual variations in continuous speech are accounted for by a method of supervised updating.

The remainder of this section discusses the definition of the reference scanning and recognition patterns, the method of supervised updating, and the application of the algorithm to continuous speech.

B. ENROLLMENT

Enrollment in the speaker-dependent word-recognition task defines the speaker-specific reference patterns for each word in the vocabulary. A total of 20 words per speaker is allowed. Each word is identified to the system and then spoken four times in isolation. These four repetitions are used to define reference scanning and recognition patterns for the word.

The enrollment strategy consists of preprocessing the data for each word, locating reference points, defining scanning patterns, and, finally, defining a recognition pattern. The preprocessing step uses the algorithm defined in Section II to provide the spectrum, energy, regression coefficients, and T-function.

The T-function is a measure of the change in the spectrum, regression coefficients and energy and for time t_i is given by:

$$\begin{split} T_{j} = & \sum_{k=1}^{2} \left[||(\overrightarrow{A}_{j+k})_{N} - (\overrightarrow{A}_{j+k-3})_{N}||^{2} + ||\overrightarrow{C}_{j+k} - \overrightarrow{C}_{j+k-3}||^{2} \right. \\ & + \frac{1}{4} ||E_{j+k} - E_{j+k-3}||^{2} + \frac{1}{4} ||E_{j+k-2} - E_{j+k-4}||^{2} \right] \end{split}$$

where

 $(\overrightarrow{A_j})_N$ = normalized amplitude vector (Appendix A) $\overrightarrow{C_j}$ = regression coefficient vector = $(c_{j1}, c_{j2})^T$ E_i = normalized scanning pattern energy. Reference points are located in each of the four enrollment words in turn for use in defining scanning patterns and recognition patterns. The steps in locating the reference points are as follows:

- (1) Locate the beginning, i_{ST}, and the end, i_{END}, of the word using an energy threshold.
- (2) Sum the energy in the word segment from is to i END to obtain SE.
- (3) Locate the time points associated with 5, 10, 90, and 95 percent of the energy sum S_E , that is $(i_5, i_{10}, i_{90}, i_{95})$.
- (4) Locate all the T-function peaks in the word segment.
- (5) If a T-function peak exists in the interval $[i_{ST}, i_{10}]$, define its location as the first reference point, RP₁; if not, let RP₁ = i_5 .
- (6) If a T-function peak exists in the interval $[i_{90}, i_{END}]$, define it as the last reference point, RP_N ; if not, let $RP_N = i_{95}$.
- (7) Generate the set, T, of all T-function peaks in the interval (RP_1, RP_N) . If T is null, then the word has only the two reference points, RP_1 and RP_N .
- (8) If T is not null, use the elements of T in all combinations to maximize the function:

$$F = \prod_{k=1}^{N} \left(\frac{T_k}{T_{min}}\right) \prod_{k=2}^{N} \left(\frac{i_k - i_{k-1}}{i_{RP_N} - i_{RP_1}}\right)^r$$

where $i_1 = RP_1$; $i_N = RP_N$; $i_k \in T$ for $k = 2 \dots, N-1$; T_k is the value of the T-function at i_k ; T_{min} is a normalization factor; and r is the power for the distance weighting. The subset of T that maximizes the function F is then used as the set of reference points for the word, with the first reference point being RP_1 and the last being RP_N . The objective of the maximization is to distribute the reference points uniformly throughout the word.

At the location of each of the reference points thus defined, a scanning pattern is defined as discussed in Section II. The scanning pattern uses the spectrum, energy and regression coefficients, and their respective differences between time samples.

The definition of the recognition patterns also makes use of the location of the reference points. The steps in defining the recognition pattern format are as follows:

- (1) If the energy is greater than a threshold for time samples $i_{RP_1} 4$, $i_{RP_1} 3$, $i_{RP_1} 2$, and $i_{RP_1} 1$, then extrapolation columns are defined at $i_{RP_1} 4$ and $i_{RP_1} 2$.
- (2) If the energy is greater than a threshold for time samples $i_{RP_N} + 1$, $i_{RP_N} + 2$, $i_{RP_N} + 3$, and $i_{RP_N} + 4$, then extrapolation columns are defined at $i_{RP_N} + 2$ and $i_{RP_N} + 4$.
- (3) Interior to every pair of reference points, i_k and i_{k+1} , interpolate with M columns for the recognition pattern, where

$$M = \begin{cases} 2 & i_{k+1} - i_k < 4 \\ 3 & 4 \le i_{k+1} - i_k \le 10 \\ 5 & 10 < i_{k+1} - i_k \le 20 \\ 8 & 20 < i_{k+1} - i_k \end{cases}$$

The format just described is used to define a recognition pattern using the procedure outlined in Section II.

Once the reference points, scanning patterns, and recognition pattern have been obtained for one of the four repetitions of a word as described above, those scanning patterns are used to scan the remaining three repetitions to automatically find reference points and define reference patterns. At each repetition, the new scanning and recognition patterns are averaged with the old patterns. Each of the four repetitions of the word are used in turn to define a set of reference scanning and recognition points. As each reference pattern is formed, a composite error consisting of the scanning error and the recognition error is computed. The minimum composite error over the four different enrollment trials defines the ultimate enrollment for the word.

An example of automatic enrollment is shown in Figure 11 for the word "Two". For an energy threshold of 100, the beginning of the word, i_{ST} , is at 26 and the end, i_{END} , is at 66. The 5-, 10-, 90-, and 95-percent energy sum points are at 29, 33, 55, and 58, respectively. Therefore, the strategy outlined above locates the reference points at 26, 33, and 58. The recognition format consists of three interpolation points between the reference points at 26 and 33, eight interpolation points between the reference points at 33 and 58, and the extrapolation beyond the reference point at 58 by two and four points.

C. UPDATING

To accommodate intersession variations and continuous-speech context variation, several sessions of supervised updating should be performed. The updating should consist of five sessions separated by at least a day. A series of phrases that contain all the transitions for the 20-word vocabulary should be spoken continuously. If the phrase is recognized, the reference patterns are updated by adding 1/16 of the new pattern to 15/16 of the old patterns. The five sessions spaced at 1-day intervals adapts the reference pattern for intersession variations. The series of phrases with different contexts allows the reference patterns to adapt to continuous speech by allowing the patterns to "see" something besides silence between the words and also to account for coarticulation which occurs in some contexts.

D. APPLICATION TO CONTINUOUS SPEECH

The method of word recognition using spectral pattern matching offers a dramatic improvement in performance compared with schemes that rely on finding word boundaries with energy profiles. The spectral pattern-matching method works well in continuous speech provided the words are enrolled properly and several sessions of continuous speech updating are accomplished, as discussed in Subsection III.C. The example in Figure 11 of a good enrollment for the word "two" shows the reference points. In the example, the registration points were interior to the word, so that the scanning patterns will not be confused with the scanning

patterns of adjacent words. With several sessions of continuous updating using phrases containing all the word transitions, the scanning and recognition patterns adapt to "see" all these transitions.

Unfortunately, the method of automatic enrollment described above does not always give a good enrollment. As an example of a poor enrollent, Figure 12 shows the automatic enrollment for the word "six". For an energy threshold of 100, the beginning of the word, i_{ST}, is at 25 and the end, i_{END}, is at 44. The 5-, 10-, 90-, and 95-percent energy sums are shown at 30, 31, 39, and 41, respectively. The automatic enrollment scheme chose the reference points at 29 and 44. The recognition pattern consists of five columns between the two reference points. There should have been another reference point at 56 and extrapolation of the recognition pattern both before the first reference point at 29 and after the last reference point at 56. As the patterns exist with the automatic enrollment, the updating will not improve the recognition of the word "six".

It is believed that an improved automatic enrollment algorithm would consist of a set of speaker-independent reference phoneme patterns. Given the phonetic spelling, the specific phoneme patterns for the word being enrolled would be scanned across the input data for that word and scanning errors obtained. The minimum scanning errors would be located and used in a dynamic programming algorithm to obtain the best sequence of phonemes in the proper spelling order for the word. T-function peaks between the minimum error locations for the phoneme pairs would be used to define the reference points for the word. The recognition pattern format would be specific to the phonemic spelling of the word. Once the reference points and recognition pattern format are defined, the enrollment procedure would be the same as defined above.

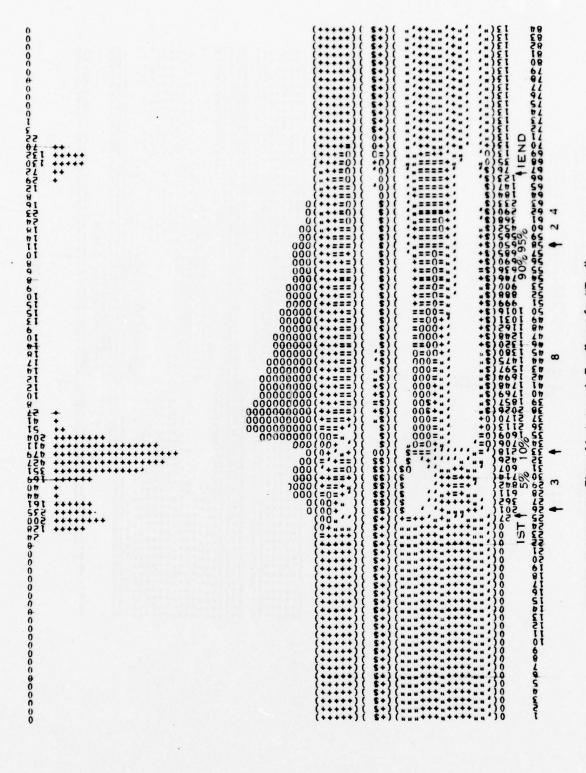


Figure 11. Automatic Enrollment for "Two"

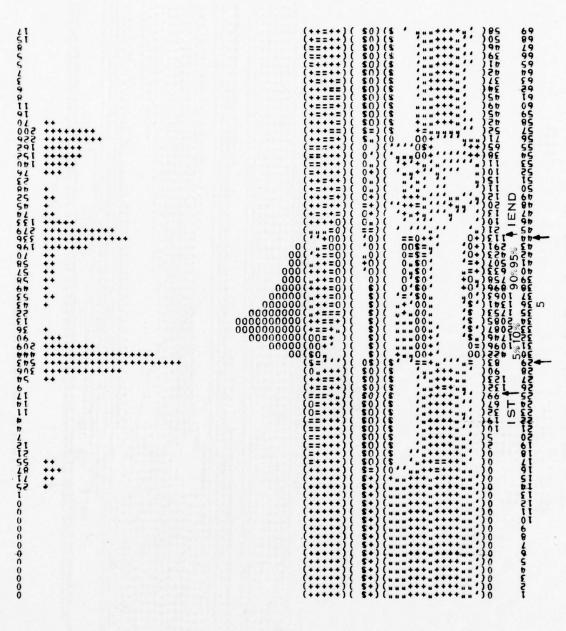


Figure 12. Automatic Enrollment for "Six"

SECTION IV

REFERENCE-PATTERN GENERATION FOR SPEAKER-INDEPENDENT WORD RECOGNITION

The creation of reference patterns for speaker-dependent, isolated word recognition is fairly straightforward: extract patterns from a single enrollment session and accommodate intersession variation through updating of the reference patterns (learning with a teacher). For speaker-independent word recognition, however, the increased variance of input data from singe-reference templates because of dialect, idiolect, and actual physical characteristics of the speaker (length and shape of the vocal tract, pitch, etc.) requires a more complex approach. Of course, allowing continuous speech input exacerbates the problem with the introduction of contextual variations. An obvious solution is to allow multiple reference templates for each word. One approach to deriving a set of multiple reference templates from a design data set is to partition the data set on the basis of information other than the actual data, such as sex or linguistic background. A second approach, the one used in the total voice verification study and in the present study, is to partition the data on the basis of the data points themselves using clustering techniques. The remainder of this section reviews the use of clustering in speakerindependent reference template generation, discusses the clustering used in the studies at Texas Instruments, and gives results of some further analysis of the patterns developed during the total voice work, patterns that were also used on the current unconstrained digit recognition work in order to preserve compatibility.

A. REVIEW OF CLUSTERING IN SPEAKER-INDEPENDENT WORD RECOGNITION

Except for the work done in this study and in the total voice study,¹ the only other applications of clustering to speaker-independent reference-pattern generation has been concurrent work started independently about the same time at Bell Laboratories (as reported by Rabiner, Levinson, Rosenburg, and Wilpon^{19 - 22}) and subsequent independent work done in Japan by Tanaka.^{23,24} The application at Bell Laboratories is isolated word recognition, and Tanaka's procedure has been applied only to the recognition of stop consonants. The remainder of this subsection reviews these other works.

¹⁹ L.R. Rabiner, "On Creating Reference Templates for Speaker-Independent Recognition of Isolated Words," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-26:34–42, February 1978.

²⁰ S.E. Levinson et al., "Interactive Clustering Techniques for Selecting Speaker-Independent Reference Templates for Isolated Word Recognition," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-27:134–141, April 1979.

²¹ L.R. Rabiner et al., "Speaker-Independent Recognition of Isolated Words Using Clustering Techniques," Proceedings of the International Conference on Acoustics, Speech and Signal Processing, Washington, D.C., 574-577, 2-4 April 1979.

^{574-577, 2-4} April, 1979.

²²L.R. Rabiner and J.G. Wilpon, "Considerations in Applying Clustering Techniques to Speaker-Independent Word Recognition," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, Washington, D.C., 578-581, 2-4 April, 1979.

²³K. Tanaka, "A Standard Category Pattern-Making Method With Application to Phoneme Recognition," Proceedings of the Fourth International Joint Conference on Pattern Recognition, Kyoto, Japan, 1030–1032, 7-10 November, 1978.

²⁴K. Tanaka, "A Talker Clustering Method for Standard Pattern Making," *Progress Report on Speech Research* '77, Electrotechnical Laboratory, Japan, August 1978.

Although the Bell Laboratories work has been for words said in isolation, the word sets investigated, while including the digits (0 through 9), were considerably larger. One set was a 54-word vocabulary proposed originally by Gold, 25 and the other set contained the alphabet, the digits, and three control words. The speech representation chosen was a set of linear predictive coding (LPC) parameters for each 15-ms frame of speech. These parameters then underwent a time normalization using a dynamic programming technique.26 The similarity measure used in the Bell Laboratories work was one of the following form proposed by Itakura:7

$$d[k, w(k)] = \log \frac{\overrightarrow{a}_{w(k)} V \overrightarrow{a}_{w(k)}}{\overrightarrow{a}_k V \overrightarrow{a}_k}$$
 (5)

where \overline{a}_k is the vector of LPC coefficients associated with the kth frame of the test or unknown utterance x_i ; $a_{w(k)}$ is the vector of LPC coefficients derived from the w(k)th frame of the reference utterance x_i; and V is the matrix of autocorrelation coefficients computed from the kth frame of the test utterance. Note that this distance measure is not a true metric since it is not symmetrical.

The clustering technique reported in Levinson et al.,20 and Rabiner et al.,21 is a supervised, interactive procedure and is the combination (figure 13) of the following four procedures: chainmap, shared nearest neighbor, k-means, and a version of ISODATA. The details of their procedures are given in Levinson et al.²⁰ In this approach, the investigators first attempted to find good estimates of both the number of clusters (using the chainmap) and their cluster centers (using the k-means) for input to an iterative optimization procedure (ISODATA) that allowed splitting and merging of clusters. The overall intent was to maximize a quality measure σ for the assignment of N observations into M classes. The value of σ is given by

$$\sigma = \frac{\frac{1}{M(M-1)} \sum_{i=1}^{M} \sum_{j=1}^{M} \delta(x_{p}^{(i)}, x_{p}^{(j)})}{\frac{1}{M} \sum_{i=1}^{M} \frac{1}{m_{i}(m_{i}-1)} \sum_{j=1}^{m_{i}} \sum_{k=1}^{m_{i}} \delta(x_{j}^{(i)}, x_{k}^{(i)})}$$
(6)

where superscripts indicate class membership, p subscripts indicate reference-class prototypes, and $\delta(a,b)$ is a nonsymmetric similarity measure between patterns a and b that is the average of the Itakura distances over all the frames of the reference pattern. Further comments regarding σ are made later in this section.

IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-26:575-582, December 1978.

²⁵B. Gold, "Word-Recognition Computer Program," Massachusetts Institute of Technology, Cambridge, RLE, Technical Report 452, June 1966.

26 L.R. Rabiner et al., "Considerations in Dynamic Time Warping Algorithms for Discrete Word Recognition,"

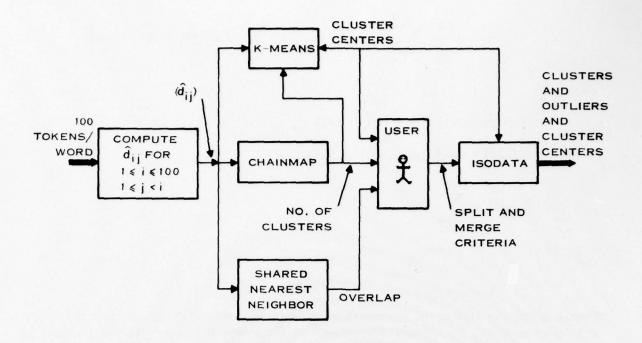


Figure 13. Bell Laboratories Clustering Procedures (From Rabiner, et al21)

Speaker-independent recognition results for isolated digits given in Rabiner et al.²¹ range from 97.5 to 100 percent. Results for the entire 39-word set range from about 50 to 80 percent, and recognition improves with the number of reference templates used for each word.

Rabiner and Wilpon²² extended the previous work to unsupervised clustering using the same data set, distance, and quality measure (σ) as previously used. One clustering algorithm uses only precomputed distances between observations, attempting to place each observation uniquely in a cluster with all others that are similar, and a second clustering algorithm combines (by averaging) observations that are similar. Comparisons were made in this work among three different LPC feature sets and between cluster representation either by the data point with the minimum maximum distance from all other points in the cluster or by the average of all the points in the cluster. The results of Rabiner and Wilpon indicate that the algorithm using precomputed distances was superior to the other and that the use of an averaged pattern to represent the cluster was superior to using the minimum maximum center. Again, the recognition accuracy improved with the number of reference templates used.

Tanaka^{23,24} clusters a set of observations into different classes by moving each observation iteratively by some amount proportional to the density of points in the neighborhood of the observation, where each point is modified by its gradient with respect to all other points. Specifically, for the set of observation vectors x_i (i = 1, ..., N) for the jth iteration,

$$\vec{x}_{i}^{j+1} = \frac{1}{w_{i}^{j}} \sum_{k=1}^{N} \delta(s^{j}) \left[\vec{x}_{k}^{j} + \frac{c}{\sqrt{w_{k}^{j}}} \sum_{r=1}^{N} \sqrt{w_{r}^{j}} \left(\vec{x}_{k}^{j} - \vec{x}_{r}^{j} \right) \delta(2s^{j}) \right]$$
(7)

where

$$\delta(s^{j}) = \exp \left\{-\left[d(x_{i}^{j}, x_{k}^{j})\right]^{2}/2(s^{j})\right\}$$

$$\delta(2s^{j}) = \exp \left\{-\left[d(x_{k}^{j}, x_{r}^{j})\right]^{2}/2(2s^{j})\right\}$$

$$w_{i}^{j+1} = \sum_{k=1}^{N} \delta(s^{j})$$

and

$$s^{j} = s^{j-1}/3\sqrt{2}$$

Tanaka makes the analogy to a potential function of an exponential form, so that the $(x^j_k - x^j_r) \delta(2s^j)$ term can be considered a gradient of the potential function. Hence, the term in brackets in Equation (7) represents one of the points modified by the gradient of a potential function of that point with respect to all other points. These modified points are then used in N weighted sums to determine each of the new N points for the (j+1)st iteration. Clustering stops when the window $\delta(s^j)$ has narrowed sufficiently that every data point is the same as that of the previous stage. The number of iterations and the final number of clusters are obviously affected by the choice of c and s^0 .

This approach is quite similar to that presented by Fukunaga and Hostetler.²⁷ However, they essentially use the points as modified in the brackets in Equation (7) as the new N points for the (j + 1)st iteration, instead of using weighted sums of such points.

Tanaka applied this method to generation of reference patterns for use in the difficult problem of detecting the three stop consonants /p,t,k/. This effort is directed toward a phonetic classification-based speech recognition system. Tanaka's results are 89 percent for the test data and 82 percent for stop consonants of other speakers.

The method used in the study covered by this final report and the method used in the preceding total voice study differ significantly from the approaches just described. Tanaka's approach differs not only in terms of using phonemic-based recognition but also in the clustering by allowing movement of the data points. Although differing final clusters and numbers of clusters could be produced in Tanaka's algorithm by varying the parameters, he does not discuss how to choose the final clusters. Differing applications do not allow comparison of his final results to those presented here.

²⁷K. Fukunaga and L.D. Hostetler, "The Estimation of the Gradient of a Density Function, With Applications in Pattern Recognition," *IEEE Transactions on Information Theory*, IT-21:32–40, January 1975.

Although the application and the clustering approach used at Bell Laboratories differ less, their approach is to find a single "best" set of clusters, ultimately using an ISODATA algorithm that can split and merge clusters. The approach used at Texas Instruments, however, finds good estimates for cluster centers using a hierarchical clustering approach, performing an iterative optimization on cluster definitions for several fixed values of M (later referred to as "c"), and choosing M based not only on criterion values for each of the final partitions, but also on a subjective evaluation of the final cluster averaged patterns themselves.

A criterion similar to that used by Bell Laboratories (σ) is compared with the one used here (trace of the within-class scatter matrix) later in this section. In addition, although a common data base could not be used, a small test of recognition performance on isolated digits was performed to provide a rough comparison with the isolated digit results presented by Rabiner et al.²¹

B. DETAILED CLUSTERING ALGORITHM

The clustering algorithm used in the total voice study and extended in the present study represents a unique combination of several methods, all centered on the use of Euclidian distances because the fast vector comparator exists peripheral to the TI 980 to perform the computation. The entire procedure is shown in Figure 14. The patterns used in the speaker-independent digit-recognition evaluations were generated during the total voice study using the path through the procedure denoted by the double lines in Figure 14. The other paths in Figure 14 were added during this study for evaluation of and consistency checks on the previous patterns and for rudimentary outlier analysis.

A detailed description of the procedure used to generate the patterns and the patterns themselves are given in the total voice final report. A brief description of the procedure is given here for completeness. The method used in the total voice study to derive the patterns used in the evaluation was an agglomerative method combining the two clusters that have the smallest average distance (MINAVE) between the points in the two clusters, i.e., combining the i and j clusters that have the minimum

$$\frac{1}{n_i n_j} \sum_{\overline{x} \in x_i} \sum_{\overline{x}' \in x_i} d(\overline{x}, \overline{x}')$$
 (8)

where

 n_i = number of \vec{x} in class χ_i n_i = number of \vec{x}' in class χ_i

and, in this case,

$$d(\vec{x}, \vec{x}') = ||x - x'||^2$$

The second step used was to improve on the partitions from the hierarchical clustering iteratively by moving samples from one group to another if such a move improved the value of

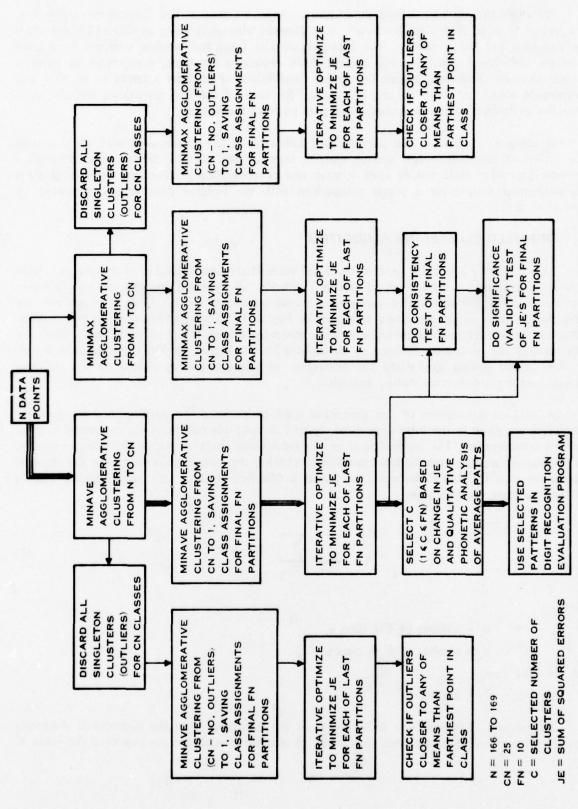


Figure 14. Block Diagram of Clustering Procedure Developed for Clustering Scanning and Recognition Patterns for Speaker-Independent Digit Recognition

some criterion function. This step used the iterative optimization method of Duda and Hart²⁸ that minimized the sum-of-squared error criterion J_e , written as

$$J_{e} = \sum_{i=1}^{c} J_{i} = \sum_{i=1}^{c} \sum_{\vec{x} \in x_{i}} || \vec{x} - \vec{m}_{i} ||^{2}$$
 (9)

where

$$\vec{m}_i = \frac{1}{n_i} \sum_{\vec{x} \in x_i} \vec{x}$$

If a point $\hat{\vec{x}}$ is moved from class x_i to class x_j , the means \vec{m}_i and \vec{m}_j change to

$$\overline{m}_i^* = \overline{m}_i - \frac{\widehat{x} - \overline{m}_i}{n_i - 1} \text{ and } \overline{m}_j^* = \overline{m}_j + \frac{\widehat{x} - \overline{m}_j}{n_i + 1}$$
 (10)

The value of J_i decreases to

$$J_{i}^{*} = J_{i} - \frac{n_{i}}{n_{i} - 1} \|\widehat{\vec{x}} - \overline{m}_{i}\|^{2}$$
 (11)

and J, increases to

$$J_{j}^{*} = J_{j} + \frac{n_{j}}{n_{j} + 1} \|\widehat{\vec{x}} - \overline{\vec{m}}_{j}\|^{2}$$
 (12)

Clearly, then, since the criterion is to minimize J_e, if

$$\frac{n_{j}}{n_{i}+1} \|\widehat{\vec{x}} - \widehat{\vec{m}}_{j}\|^{2} < \frac{n_{i}}{n_{i}-1} \|\widehat{\vec{x}} - \widehat{\vec{m}}_{i}\|^{2}$$
 (13)

then \hat{x} should be transferred from class x_i to class x_j . Specifically the point \hat{x} is moved to the class x_j , having the smallest $(n_i/n_i + 1) ||\hat{x} - m_i||^2$.

An additional property (not necessarily good) of the selection of J_e as a criterion is that a set of equally divided clusters is favored over a set containing both small and large clusters, as noted earlier. This can be seen by considering $n_i \gg n_j$ in Equation (13), which yields approximately

$$\frac{n_{j}}{n_{i}+1} \|\widehat{\vec{x}} - \vec{m}_{j}\|^{2} < \|\widehat{\vec{x}} - \vec{m}_{i}\|^{2}$$
 (14)

²⁸ R.O. Duda and P.E. Hart, Pattern Classification and Scene Analysis, John Wiley and Sons (New York, 1973).

Thus, for $n_j = 1$, the distance $||\widehat{x} - m_j||^2$ need only be less than twice the distance $||\widehat{x} - m_i||^2$ to the old mean for \widehat{x} to the transferred to class x_j .

C. CRITERIA FOR MEASURING PARTITION GOODNESS

Although minimization of J_e was the criterion used in the iterative optimization and $(J_e^c - J_e^{c+1})/J_e^c$ was used as a second criterion to aid in selecting the number of clusters, the values of several other criteria related to J_e were calculated during the current study for all patterns for numbers of classes from 1 to FN (= 10). (Superscripts on J_e are used to denote number of classes.)

The discussion in the remainder of this subsection assumes a knowledge of scatter matrices. Appendix B has been provided for those not familiar with the concept.

A third criterion is the value of tr $S_B/\text{tr }S_W$, which is inherently maximized during the iterative optimization by the minimization of J_e (= tr S_W). Note that since tr $S_T = J_e^{-1}$ and tr $S_B = \text{tr }S_T - \text{tr }S_W$, then tr $S_B/\text{tr }S_W = (J_e^{-1} - J_e^c)/J_e^c$ for c classes.

A fourth related criterion suggested by Hartigan²⁹ for choosing the number of clusters is $(n-c)(J_e^c-J_c^{c+1})/J_e^c$. Hartigan suggests that values of this ratio greater than 10 justify increasing the number of clusters.

A fifth criterion is related to the F-ratio from analysis of variance, taking into account the degrees of freedom of tr S_B and tr S_W . This criterion is given by

$$\frac{\operatorname{tr} S_{B}/(c-1)}{\operatorname{tr} S_{W}/(n-c)} = \frac{(n-c)(J_{e}^{-1} - J_{e}^{c})}{(c-1)J_{e}^{c}}$$
(15)

and is attribted by Everitt³⁰ to Calinski and Harabasz.³¹

The sixth criterion calculated during this study is a σ analogous to that used in the Bell Laboratories studies.¹⁹⁻²² The value of σ is calculated by

$$\sigma = \frac{\frac{1}{c(c-1)} \sum_{i=1}^{c} \sum_{j=1}^{c} \|\vec{m}_{i} - \vec{m}_{j}\|^{2}}{\frac{1}{c} \sum_{i=1}^{c} \frac{1}{(n_{i}-1)} \left[\frac{1}{n_{i}} \sum_{\vec{x} \in X_{i}} \sum_{\vec{x}' \in X_{i}} \|\vec{x} - \vec{x}'\|^{2} \right]}$$
(16)

²⁹ J.A. Hartigan, Clustering Algorithms, John Wiley and Sons (New York, 1975).

³⁰ B. Everitt, Cluster Analysis, Heinemann Educational Books, Ltd., (London, 1974).

³¹ T. Calinski and J. Harabasz, "A Dendrite Method for Cluster Analysis," (unpublished), 1971.

The relationship between σ and the tr S_B/tr S_W criterion can be seen better by putting both σ and criterion 5 in equivalent forms. The σ term can be rewritten as

$$\sigma = \frac{\frac{c}{c-1} \sum_{i=1}^{c} \sum_{j=1}^{c} \frac{1}{c^{2}} \|\vec{m}_{i} - \vec{m}_{j}\|^{2}}{\frac{2}{c} \sum_{i=1}^{c} \frac{1}{(n_{i}-1)} J_{i}}$$
(17)

and the fifth criterion multiplied by the factor c/(2n) can be rewritten (see Appendix C) as a seventh criterion α , as follows:

$$\alpha = \frac{c(\text{tr } S_B)/(c-1)}{2n(\text{tr } S_W)/(n-c)} = \frac{\frac{c}{c-1} \sum_{i=1}^{c} \sum_{j=1}^{c} \frac{n_i n_j}{n^2} \|\vec{m}_i - \vec{m}_j\|^2}{\frac{2}{c} \sum_{i=1}^{c} \frac{1}{\binom{n}{c}-1} J_i}$$
(18)

The values of all seven of these criteria for several classes of one of the 34 pattern types clustered in this study are shown in Table 5. The desire is to maximize the last five of the seven criteria discussed above. Note, however, that the sixth criterion, σ , is actually only monitored, while the optimization is on the basis of J_e , from which all the other criteria (except σ) are derived.

TABLE 5. VALUES OF SEVEN CLUSTERING CRITERIA FOR POSTITERATIVE OPTIMIZATION OF MINMAX AGGLOMERATIVE CLUSTERS FOR RECOGNITION PATTERN FOR DIGIT "SIX"

	Criterion									
<u>c</u>	1	2	3	4	5	6	7			
1	58,412.0	0.176	0.000	28.850	0.000	0.000	0.000			
2	48,136.5	0.051	0.213	8.296	34.795	0.422	0.208			
3	45,686.7	0.055	0.279	8.897	22.561	0.426	0.203			
4	43,177.7	0.022	0.353	3.473	18.935	0.518	0.227			
5	42,246.4	0.030	0.383	4.771	15.306	0.530	0.229			
6	40,986.8	0.019	0.425	2.943	13.520	0.486	0.243			
7	40,228.1	0.036	0.452	5.652	11.903	0.546	0.249			
8	38,789.0	0.003	0.506	0.543	11.346	0.601	0.272			
9	38,654.9	0.021	0.511	3.343	9.987	0.601	0.269			
10	37,826.5	-	0.544		9.372	0.633	0.281			

D. DESCRIPTION OF CLUSTER ANALYSIS DOCUMENTATION

This subsection gives a brief description of the printouts available from an analysis program that was run on the output data from the entire clustering procedure shown in Figure 9 for each of the 34 pattern types. A subset of the available outputs is presented in Appendix D. Samples of available printouts are given in this subsection for the scanning pattern for reference-point 1 for the digit zero. MINAVE is used to refer to agglomerative clustering by combining the two clusters having the minimum average distance between all points in the two clusters. Corespondingly, MINMAX refers to agglomerative clustering by combining the two clusters having the minimum maximum distance between the two clusters.

1. Trees

The first type of output available is a tree (dendogram) showing the final joinings or agglomerations in the hierarchical procedures and is available for all four branches of the algorithm shown in Figure 9. Accompanying each dendogram is a table showing the values of the joining criterion for each level and the relationship of the criterion values to the dendogram. The tree-printing subroutine was adapted from appendix G of Anderberg.³² The tree for the MINAVE hierarchical clustering using all samples is shown in Figure 15 for the joining criteria values in Table 6.

2. Parameter Comparisons

The second type of output from the analysis program gives the values of the six criteria described in Subsection IV.C, the values of the errors during the agglomerative clustering, and the number of iterations required in the iterative optimization to reach the final partitions. The conditions for each of the parameter comparisons produced and a reference to the figure showing an example of that comparison are listed below:

MINAVE, all points; pre- to postiterative optimization comparison (Figure 16)
MINAVE, outliers discarded; pre- to postiterative comparison (Figure 17)
MINMAX, all points; pre- to postiterative optimization comparison (Figure 18)
MINMAX, outliers discarded; pre- to postiterative comparison (Figure 19)
Preiterative optimization, all points; MINAVE to MINMAX comparison (Figure 20)
Postiterative optimization, all points; MINAVE to MINMAX comparison (Figure 21).

3. Consistency Tests

The class assignments obtained using the MINAVE and the MINMAX hierarchical clustering procedures after iterative optimization are compared for the number of clusters ranging from 2 to FN (= 10). This comparison is in terms of two contingency matrices such as shown in Figure 22 for 10 classes after iterative optimization. This output first lists the members of each class for the iteratively optimized results of the MINAVE agglomerative clustering followed by those from the MINMAX agglomerative clustering. The first contingency table then compares the

³²M.R. Anderberg, Cluster Analysis for Applications, Academic Press (New York, 1973).

TABLE 6. CRITERION VALUES FOR TREE FOR FINAL 24 STAGES OF MINAVE AGGLOMERATIVE CLUSTERING OF ALL (166) SCANNING PATTERNS FOR FIRST REFERENCE POINT OF DIGIT "ZERO"

	Cl	ass	Criterion		
Stage	1	J	Absolute	Relative	
1	16	22	403.692	1	
2	7	23	406.000	1	
3	2	6	408.710	1	
4	7	15	414.500	1	
5	10	19	414.750	1	
6	9	21	418.000	2	
7	1	11	419.967	2	
8	9	18	429.875	3	
9	13	20	439.000	3	
10	2	3	441.115	3	
11	4	8	461.250	5	
12	14	16	462.714	5	
13	1	12	464.742	5	
14	1	2	488.888	7	
15	1	5	509.195	9	
16	1	10	527.124	10	
17	4	24	530.000	10	
18	1	14	538.098	11	
19	4	25	542.947	11	
20	1	13	573.328	14	
21	1	9	626.332	18	
22	1	7	655.048	20	
23	4	17	664.150	20	
24	1	4	729.398	25	

members of the classes from each partition, with each entry in the table showing the number of points that are members of both classes. The second contingency table compares data points in pairs for joint membership or lack of joint membership in the same class. In particular, if two data items are in different classes in a partition, this fact is denoted by a 1 in row or column 1. Otherwise, a 1 appears in row or column 2. Hence, the (2,2) entry in the contingency table indicates how many pairs of the N(N-1)/2 pairs of points are in the same class for both partitions, and the (1,1) entry indicates how many of the pairs are in different classes for both partitions.

Ideally, for both contingency tables, all off-diagonal elements will be 0. Hence, a measure of the closeness of the two partitions in both cases is the sum of the diagonal entries divided by the sum of all the entries in the table [N for the first table and N(N-1)/2 for the second table].

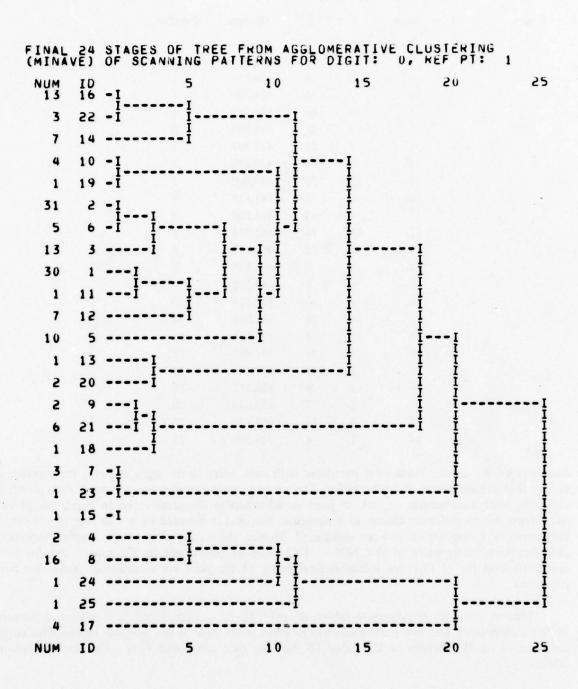


Figure 15. Tree for Final 24 Stages of MINAVE Agglomerative Clustering of All (166) Scanning Patterns for First Reference Point of Digit "Zero"

STATISTICS FOR DIGIT: 0; REF PT: 1; NO OF DATA PIS: 166
MINAVE AGGLOM CLUSTERING; 23 FEB79
PRE AND POST ITERATIVE OPTIMIZATION FOR MIN JE

С	ERR ITERS	JE (=TR(w))	JE(C)-JE(C+1) /JE(C)	IK(R)\15(M)
1234567890123456789012345	0.0 729.4 80 939 1324 141 96 655.3 141 96 655.3 141 96 97 141 96 97 141 96 97 141 96 97 141 141 141 141 141 141 141 14	PRE 40504.9 40730.9 40504.6 284.8 40730.9 39361.8 333624.6 330504.6 330604.6 330604.6 330604.6 330604.6 330604.6 330604.6 330604.	PRE	PRE 0.000 0.000 0.143 0.176 0.278 0.417 0.578 0.599 0.517 0.313 0.599 0.402 0.644 0.422 0.690 0.720 0.518 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.744 -0.000 0.745 -0.000 0.953 -0.000 0.973 -0.000 0.000 0.973 -0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.000 0.
c	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C) *TR(B) /(C-1) *TR(W)
123456789012345	0.000 0.000 0.135 0.141 0.139 0.2481 0.139 0.2481 0.173 0.299 0.175 0.3350 0.224 0.3567 0.2239 0.366-0.000 0.371-0.000 0.371-0.000 0.371-0.000 0.459-0.000 0.459-0.000 0.466-0.000 0.459-0.000 0.466-0.000 0.467-0.000 0.487-0.000 0.5157-0.000	PRE 0.000 0.000 0.000 0.579 0.499 1.042 0.614 1.114 0.6644 1.1153 0.6646 1.409 0.7727 1.585 1.533***********************************	PRE 19.799 0.608 19.799 20.608 55.512 19.406 55.147 9.310 8.083 7.330 10.460 5.201 1.790 4.316 3.815 2.628 7.330 10.460 7.300 1.790 4.712 0.000 4.712 0.000 1.566 0.000 1.566 0.000 1.566 0.000 1.566 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000 1.568 0.000	PRE 0.000 0.

Figure 16. Parameter Comparisons for Pre- and Postiterative Optimization of MINAVE Partitions Using All Points

STATISTICS FOR DIGIT: U; REF PT: 1; NO OF DATA PTS: 158
MINAVE AGGLOM CLUSTERING; 23 FEB79
PRE AND POST ITERATIVE OPTIMIZATION FOR MIN JE

C	ERR ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)
123456789011234567	0.0 718.9 681.5 681.5 79 651.6 101 562.4 653 536.4 74 477.1 62 477.1 655 461.3 73 451.4 447.0 435.2 418.0 -1 408.7 -1 408.7 -1	PRE 146284 9 37601 3 37357 9 36595 8 32946 4 4 35214 9 30736 0 0 32448 7 286913 - 9 28697 2 26819 1 24448 6 26	PRE	PRE 0.000 0.000 0.142 0.239 0.173 0.405 0.501 0.
С	C(N-C)TR(B) /2N(C-1)TR(N)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) * TR(B) /(C-1) * TR(W)
123456789011234567	0.000 0.000 0.133 0.225 0.121 0.284 0.135 0.313 0.186 0.370 0.263 0.384 0.280 0.406 0.303 0.451 0.311 0.474 0.317 0.000 0.408 - 0.000 0.408 - 0.000 0.424 - 0.000 0.434 - 0.000 0.456 - 0.000	PRE 0.000 0.000 0.000 0.630 0.552 1.040 0.664 1.371 0.635 1.204 0.733 1.204 0.733 1.130 0.976 1.150 0.921 1.173 1.203************************************	PRE PUST 19.472 30.281 4.172 18.422 5.849 10.399 12.097 9.130 1.690 6.452 4.366 4.534 5.336 5.478 8.483 2.488 4.879 2.034****** 12.526 0.000 2.424 0.000 2.424 0.000 2.528 0.000	PRE 0.000 0.000 22.087 37.277 13.404 31.376 11.2349 22.969 13.142 20.474 12.477 18.195 11.616 16.845 11.185 16.635 10.322 15.741 9.561 -0.000 10.451 -0.000 9.615 -0.000 9.436 -0.000 9.436 -0.000 9.436 -0.000 9.103

Figure 17. Parameter Comparisons for Pre- and Postiterative Optimization of MINAVE Partitions With Outliers Discarded

STATISTICS FOR DIGIT: 0; REF PT: 1; NO OF DATA PTS: 166
MINMAX AGGLOM CLUSTERING; 23 FEB79
PRE AND POST ITERATIVE OPTIMIZATION FOR MIN JE

С	ERR ITERS	JE (=TR(w))	JE(C)-JE(C+1) /JE(C)	TK(B)/TR(W)
1234567890123456789012345	0.0 1421.0 54 1331.0 23 1091.0 76 1020.0 84 976.0 114 862.0 8844.0 805.0 804.0 -1 756.0 -1 712.0 -1 682.0 -1 682.0 -1 682.0 -1 682.0 -1 682.0 -1 683.0 -1 683.0 -1 684.0 -1 685.0 -1 685.0 -1 686.0 -1 6	PRE 401344.8 9 42816.1 337736.3 337736.3 35985.1 337745.3 3736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 373736.3 3737373737373737373737373737373737373	PRE	PRE
C	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) *TR(B) /(C-1) *TR(W)
12345678901123456789012345	0.000 0.000 0.080 0.151 0.188 0.217 0.186 0.248 0.214 0.304 0.271 0.310 0.287 0.344 0.271 0.390 0.309 0.370 0.354-0.000 0.405-0.000 0.405-0.000 0.405-0.000 0.405-0.000 0.405-0.000 0.405-0.000 0.405-0.000 0.517-0.000 0.557-0.000 0.557-0.000 0.5595-0.000 0.5695-0.000 0.615-0.000	PRE 0.000 0.000 0.000 0.167 0.000 0.6527 0.6525 0.6623 0.6623 0.6623 0.6633 0.6	PRE 1.924 1.924 1.924 1.924 1.924 1.924 1.924 1.924 1.924 1.924 1.926 1.924 1.926 1.924 1.926 1.	PRE 0.000 0.000 13.286 24.063 15.455 20.086 14.264 14.264 11.3330 14.264 11.393 13.656 10.685 -0.000 9.588 -0.000 9.588 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.372 -0.000 9.373 -0.000 9.374 -0

Figure 18. Parameter Comparisons for Pre- and Postiterative Optimization of MINMAX Partitions Using All Points

STATISTICS FOR DIGIT: 0; REF PT: 1; NO OF DATA PTS: 164
PRE AND POST ITERATIVE OPTIMIZATION FOR MIN JE

С	ERR ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)
12345678901234567890123	0.0 0 1421.0 53 1331.0 24 1020.0 49 9762.0 43 857.0 49 8555.0 75 8444.0 60 8055.0 -1 712.0 -1 682.0 -1 682.0 -1 678.0 -1 678.0 -1 678.0 -1 577.0 -1 577.0 -1 5571.0 -1	PRE	PRE	PRE
C	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C)*TR(B) /(C-1)*TR(W)
12345678901234567890123	0.000 0.000 0.082 0.174 0.184 0.235 0.204 0.268 0.227 0.319 0.254 0.317 0.282 0.401 0.341 0.423 0.359-0.000 0.406-0.000 0.418-0.000 0.446-0.000 0.498-0.000 0.498-0.000 0.517-0.000 0.5575-0.000 0.575-0.000 0.575-0.000 0.575-0.000	PRE POST 0.000 0.172 0.305 0.573 0.6420 0.5774 0.6449 0.574 0.6476 0.820 0.832 0.899 0.899 0.911 0.8897 0.985 ************************************	PRE 24.702 21.917 17.684 7.920 10.618 7.145 8.359 7.234 6.730 6.897 6.730 5.467 5.495 4.483 4.878 4.241****** 6.289 0.000 3.5720 0.000 3.5720 0.000 3.5720 0.000 4.140 0.000 4.140 0.000 4.140 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000 4.195 0.000	PRE 0.000 0.000 13.579 26.935 26.005 16.973 22.211 15.100 17.657 13.387 16.436 17.657 12.417 15.619 11.891 14.782 11.324 14.035 10.845 -0.000 9.845 -0.000 9.626 -0.000 9.719 -0.000 9.544 -0.000 9.719 -0.000 9.544 -0.000 9.719 -0.000 9.544 -0.000 9.719 -0.000 9.544 -0.000 9.719 -0.000 9.544 -0.000 9.719 -0.000 9.544

Figure 19. Parameter Comparisons for Pre- and Postiterative Optimization of MINMAX Partitions With Outliers Discarded

STATISTICS FOR DIGIT: 0; REF PT: 1; NO OF DATA PTS: 166
MINAVE AND MINMAX AGGLOM CLUSTERING; 23 FEB79
PRE ITERATIVE OPTIMIZATION FOR MIN JE

С	ERR	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)		
1234567890123456789012345	A VE	39361.8 35928745.3592878.3592878.3 321388.3 321388.2 3288797.3 3287359.9 7 2265544.3 321544.3	AVE	AVE		
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C)*TR(B) /(C-1)*TR(W)		
1234567890123456789012345	AVE 0.000 0.133 0.138 0.138 0.128 0.138 0.166 0.228 0.173 0.226 0.173 0.226 0.336 0.348 0.348 0.348 0.348 0.348 0.348 0.348 0.348 0.348 0.348 0.4418 0.4473 0.448 0.4473 0.4470 0.4470 0.4470 0.4470 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4480 0.4502 0.5502 0.5502 0.5502 0.5502 0.5502 0.5502 0.5502	A VE 0.000 0.579 1.042 1.153 0.6991 1.153 0.6991 1.161 0.6991 1.153 0.6991 1.16585 0.7999 1.327 0.6884 1.327 0.6884 1.327 1.6585 1.061 1.327 1.6585 1.0888 1.1538 1.16585 1.1688 1.1693 1.1688 1.1693 1.1688 1.1693 1.1688 1.1693 1.1688 1.1693 1.	A VE 9 1 22 48 23 2 1 4 1 4 2 4 5 6 1 1 2 2 3 8 9 1 3 2 1 4 1 4 7 6 6 1 1 4 1 4 2 1 4 1 2 1 4 1 4 1 4 1 4 1 4	A VE 0.000 13.1243 14.158 15.2643 11.4250 14.0164 11.050		

Figure 20. Parameter Comparisons for Preiterative Optimization of MINAVE and MINMAX Partitions Using All Points

STATISTICS FOR DIGIT: 0; REF PT: 1; NO OF DATA PTS: 166
MINAVE AND MINMAX AGGLOM CLUSTERING; 23 FEB79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO OF ITERS	JE (=TR(w))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)	
12345 6789 10	AVE MAX 0 0 8 54 80 23 93 76 139 84 124 114 141 44 96 68 80 60 75 48	AVE MAX 46284.9 46284. 40504.1 40134. 35711.4 35734. 33671.6 33736. 316504.0 30345. 29106.6 29771. 28154.5 28361. 27385.4 27293. 26902.9 26491.	8 0.118 0.110 2 0.057 0.056 4 0.061 0.060 3 0.036 0.044 3 0.046 0.019 7 0.033 0.047 9 0.027 0.038 4 0.018 0.029	AVE MAX 0.000 0.000 0.143 0.153 0.296 0.295 0.375 0.372 0.463 0.459 0.517 0.525 0.590 0.555 0.644 0.632 0.690 0.696 0.720 0.747	
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C) * TR(B) /(C-1) * TR(N)	
12345 6789 10	AVE 0.000 0.139 0.150 0.215 0.215 0.215 0.275 0.277 0.275 0.295 0.300 0.326 0.306 0.346 0.339 0.362 0.365 0.371 0.385	AVE MAX 0.000 0.000 0.499 0.305 0.614 0.627 0.616 0.633 0.644 0.663 0.702 0.816 0.772 0.839 0.797 0.880 0.864 0.891	19.169 17.763 9.196 9.001 9.710 9.533 5.652 6.921 7.238 2.986 5.135 7.435 4.262 5.877	AVE 0.000 0.000 23.121 24.824 23.835 23.768 19.979 19.838 18.420 18.599 15.443 14.514 14.351 14.083 13.371 13.482 12.328 12.785	

Figure 21. Parameter Comparisons for Postiterative Optimization of MINAVE and MINMAX Partitions Using All Points

After all class assignments and contingency tables have been printed, the measure just described is printed in a summary listing for the partition comparisons both before and after optimization (Figure 23).

E. OBSERVATIONS ON THE CLUSTERING RESULTS

Since the conclusions to be presented in this subsection have not been proved analytically but have only been observed from the clustering results, they are presented as observations only. However, these observations are made on the basis of the results of clustering 34 different pattern types, which would imply some generality to the observations.

In investigating the properties of various criteria, it is first useful to have a measure of the distribution of the class size, i.e., whether most of the samples are in one class or whether they are evenly distributed in all classes. The measure used was a normalized version of the entropy given by

$$E = -(1/EMAX) \sum_{i=1}^{c} (n_i/n) \log_2 (n_i/n)$$

$$n = \sum_{i=1}^{c} n_i$$
(19)

where

```
MINAVE
154 155
128 133
19 110
110 130
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 131
130 130 131
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 130 130
130 
                                                                                                                                                                      17
74
                                                                                                                                                                                                                                       250
149
148
148
148
155
157
162
163
143
                                                                                                                                                                                                                                                                                                                                                                                                                                                                      106
                              1
                                                                                                                                                             1744
1444
11186
11186
11188
1487
1166
109
                                                                                                                                                                                                   148
124
141
143
150
151
152
7
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  98
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      99
                              2
                                                                                                                                                                                                                                                                                                                                                                                                                                            88
                              3
                                                                                                                                                                                                                                                                                        37
                                                                                                                                                                                                                                                                                                                            45
                                                                                                                                                                                                                                                                                                                                                                                                        96
                                                                                                                                                                                                                                                                                                                                                                 61
                                                                                                                                                                                                                                                                                                                                                                                                                                            97
                                                                                                                                                                                                                                                                                                                                                                                                                                                                      114
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            115
                             45
                                                                                                                                                                                                                                                                             139
70
164
137
63
163
134
146
                                                                                                                                                                                                                                                                                                                                                               76
                                                                                                                                                                                                                                                                                                                                                                                                       77
                                                                                                                                                                                                                                                                                                                                                                                                                                            79
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  80
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      81
                             67
                                                                                                                                                                                                                                                                                                                                                                 67
                                                                                                                                                                                                                                                                                                                                                                                                        71
                                                                                                                                                                                                                                                                                                                                                                                                                                            82
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  85
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      86
                                                                                                                                                                                                                                                                                                                            66
                                                                                                                                                                                                                                                                                                                 135
40
152
                                                                                                                                                                                                                                                                                                                                                                                           57
154
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      69
                                                                                                                                                                                                 116
                                                                              MINMAX
15 61
144 1457
155 155
1260 159
95 104
130 130
131 125
131 178
150 119
150 119
150 119
150 119
150 119
150 119
150 119
                                                                                                                                                             CLASSES
17 18
74 96
165 166
38 39
                                                                                                                                                                                                                                         100
                                                                                                                                                                                                                                                                              101
                                                                                                                                                                                                                                                                                                                                                        126
                                                                                                                                                                                                                                                                                                                  25
                                                                                                                                                                                                                                                                                                                                                                                           128
                                                                                                                                                                                                                                                                                                                                                                                                                               129
                              1
                                               1332
149
149
911
1349
                                                                                                                                                                                                                                                   75
                                                                                                                                                                                                                                                                                                                                                        102
                                                                                                                                                                                                                                                                                        90
                                                                                                                                                                                                                                                                                                                                                                                                                                                                        146
                             2
                                                                                                                                                                                                                                                                                                                              91
                                                                                                                                                                                                                                                                                                                                                                                           103
                                                                                                                                                                                                                                                                                                                                                                                                                                127
                                                                                                                                                                                                                                         143
                                                                                                                                                                                                                                                                                                                                                                  42
54
97
                                                                                                                                                                                                                                                                                                                                                                                                                                   65
160
115
                              3
                                                                                                                                                                                                                                                                                      29
52
46
                                                                                                                                                                                                                                                                                                                                                                                            57
157
114
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      64
                                                                                                                                                                                                                                                                                                                   153
60
                                                                                                                                                              109
1408
1428
147
151
151
89
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  68
                                                                                                                                                                                                  1167
141
156
107
168
168
168
1102
                                                                                                                                                                                                                                                                                                                                                                                                                                                                         122
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              132
                                                                                                                                                                                                                                                                             70
164
137
64
63
                                                                                                                                                                                                                                       59
158
117
56
62
163
121
113
                                                                                                                                                                                                                                                                                                                            73
                             5
                                                                                                                                                                                                                                                                                                                                                                   76
                                                                                                                                                                                                                                                                                                                                                                                                        77
                                                                                                                                                                                                                                                                                                                                                                                                                                            79
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  80
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      81
                                                        84
10
13
39
48
14
                                                                                                                                                                                                                                                                                                                                                                                            108
                                                                                                                                                                                                                                                                                                                                                                                                                                  110
                                                                                                                                                                                                                                                                                                                              85
                                                                                                                                                                                                                                                                                                                                                                                                                                                                       111
                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      93
                                                                                                                                                                                                                                                                                                                              66
                                                                                                                                                                                                                                                                             123
                                                                                                                                                                                                                                                                                                                  148
138
                                                                                                                                                                                                                                                                                                                                                        139
 C1:
                                                                                                        MINMAX CLASSES
                                                                                                                         3
                                                                                                                                                                                                                                                                      1500
MINAVE
                                                                            9157248603
                                                                                                                2 0000000000
                                                                                                                                                     21 00 00 00 00 00 10 0
                                                                                                                                                                                           2000000000
                                                                                                                                                                                                                                00010900100
                                                                                                                                                                                                                                                                                                                                                                                             00000007000
                                                                                                                                                                                                                                                                                                                                                        00009009
 C2:
                                                                                                        MINMAX CLASSES
                                                                                                                                                                                                                                          5
                                                                                                                                                                                                                      332
MINAVE
                                                                                                                                   11618
```

Figure 22. Class Assignments and Contingency Tables for Postiterative Optimization of 10 Classes Using All Points for Reference Point 1 of Digit 0

		PRE	ITERAT	IVE OP	ITAISA	TION FO	NIM NC	JE	
NC:	5	3	4	5	6	7	8	9	10
C1:	0.440	0.548	0.542	0.536	0.470	0.470	0.392	0.416	0.434
cs:	0.504	0.620	0.618	0.598	0.573	0.540	0.643	0.651	0.676
		POST	ITERAT	VE OP	TIMIZA	TION FO	OH MIN	JE	
NC:	5	5	4	5	6	7	8	9	10
C1:	0.633	0.964	0.657	0.795	0.699	0.627	0.633	0.651	0.735
cs:			0.736						0.937
	E:a	**** 32 Ca	atinganar	Table Man	curac for D	en and Da	etitarativa		

and EMAX = $\log_2 c$. This function is maximum for c equally divided clusters having $n_i = n/c$, and is minimum for (c-1) $n_i's = 1$ clusters having one cluster $n_i = (n-c-1)/n$.

Optimization Using All Points for Reference Point 1 of Digit 0

Observation 1: The first observation is experimental confirmation of the known fact in cluster analysis that the minimization of J_e , the sum-of-squared error, favors equal sized clusters, an example of which is given in Figure B-1. One demonstration of why this is true is given in Subsection IV.B. To demonstrate in another way, remember that minimizing J_e is equivalent to maximizing tr S_B , given by

$$\operatorname{tr} S_{B} = \frac{1}{2n} \sum_{i=1}^{c} \sum_{j=1}^{c} n_{i} n_{j} \| \overline{m}_{i} - \overline{m}_{j} \|^{2}$$
 (20)

Consider just one of the terms. Assuming that the means remain relatively constant as points are changed between class i and class j, it is easy to show that $n_i n_j$ is maximum for $n_i = n_j$.

The results of the iterative optimization based on minimizing J_e in fact showed that E (the normalized entropy) increased after the optimization as shown by the histograms of E for the MINAVE agglomerative clustering in Figure 24 for the 24 scanning patterns (c = 2 through 10) and in Figure 25 for the 10 recognition patterns (c = 2 through 10).

Observation 2: MINMAX agglomerative clustering favors equal-sized clusters more than MINAVE agglomerative clustering. No previous reference to this type of observation could be found in the clustering literature, although many references existed to the "chaining" that occurs in agglomerative clustering when combining the two clusters having the minimum minimum distance between them.

Figure 26 shows the results as histograms of E for both the MINAVE and the MINMAX agglomerative clustering (before iterative optimization) for the 24 scanning patterns (c = 2 through 10). The results shown in Figure 27 are the same for the recognition patterns.

Observation 1 indicated that lower $J_e s$ (and, hence, larger αs) result from more equal-sized clusters, so that two "corollaries" exist for observation 2. The first is that J_e for MINMAX agglomerative clustering (preiterative optimization) is generally smaller than J_e for MINAVE agglomerative clustering. The second corollary is that, since J_e is usually lower for the MINMAX agglomeration, the number of iterations is also lower for the iterative optimization of the MINMAX agglomerative clusters. Histograms are not given, but refer to the columns for the number of iterations given in the 10 tables of postiterative optimization statistics for the 10 recognition patterns in Appendix D.

Observation 3: In spite of quite different starting partitions, the iterative optimization procedure applied to the MINAVE agglomerative clusters and to the MINMAX agglomerative clusters yielded similar partitions. A measure of similarity given by the ratio of the sum of the diagonal entries of a contingency table divided by the sum of all the entries in the table is shown in Table 7 for both scanning and recognition patterns before iterative optimization and in Table 8 for both pattern types after iterative optimization. These similarity measures are for the first type of contingency table described in Subsection IV.D.3.

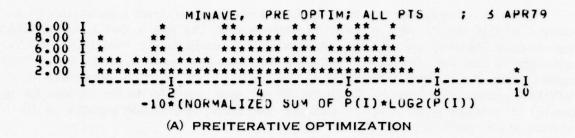
Observation 4: Although α was larger after iterative optimization than before (as it should be since the optimization criterion is to minimize J_e , which is equivalent to maximizing α), σ almost always decreased for the optimization using the MINAVE agglomerative beginning partitions and decreased for about half the pattern types for optimization using the MINMAX agglomerative beginning partitions. Plots of $\sigma_{\rm post}/\sigma_{\rm pre}$ versus $\alpha_{\rm post}/\alpha_{\rm pre}$ are given in Figure 28 for the MINAVE clustering and in Figure 29 for the MINMAX clustering for scanning patterns.

Two important points should be made. The first point is that the σ decrease was much greater for the MINAVE case than for the MINMAX case because σ_{pre} was much larger for the former, as can be seen by comparing Figure 30 and Figure 31, which show σ_{pre} versus α_{pre} for scanning patterns in both cases. As a matter of interest, Figure 32 shows σ versus α for postoptimization of the MINAVE agglomerative clustering partitions, showing that as the clusters approach equal sizes, as happens for the optimization (observation 1), σ and α approach the same value. In fact, for equally divided cluster partitions, $\sigma = \alpha$, as can be seen from the final two expressions in Subsection IV.C by setting $n_i = n_j = n/c$. The second point to remember is to temper the conclusions reached about the relationship between σ and α by the fact that optimization is with respect to α (actually J_e) while σ is being monitored only. Possibly different conclusions would be reached if iterative optimization were with respect to σ with α being monitored only.

F. TESTING CLUSTER VALIDITY WITH A PRIORI INFORMATION ABOUT DATA

The problem of testing cluster validity is a subject that has received very little attention in the literature, probably because of the difficulty of the problem. One of the few references is in Duda and Hart,²⁸ who use a hypothesis testing approach to test validity on the basis of the size of the reduction in J_e. In the specific example given, they assume multivariate normal distributions and advance the hypothesis that the data are actually from one cluster. They then derive an expression for testing this hypothesis for J_e to a specified significance level.

The approach taken in this section is an entirely different method for testing cluster validity. In the total voice speaker verification final report, descriptions of the characteristics of the reference patterns generated from the clustering algorithm are given in terms of a priori



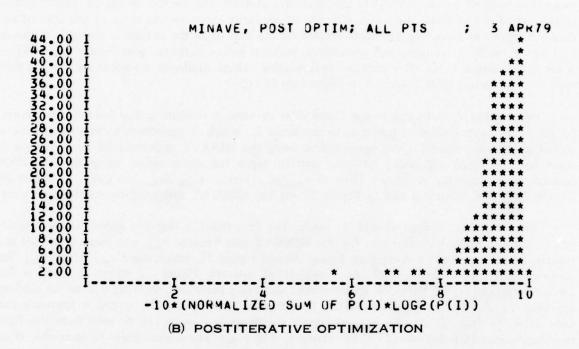


Figure 24. Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes for MINAVE Agglomerative Clustering of Scanning Patterns

information known about the data points making up each class. In this study, a quantitative measure is used for the same type of comparison. Specifically, it is assumed that a male/female division of the data is a correct way to separate the data, and then the degree to which the actual clusters agree with this assumption is measured. Because of the differences in vocal tract resonances (formants) between males and females, this is a good assumption in most cases (probably a better assumption than assuming a unimodal distribution for the data). Reference to the total voice final report, however, reveals cases where the data actually cluster on the basis of other attributes such as the scanning patterns for the third reference point of "two" which, since the formants for /u/ for males and females are very close, splits according to context. This would suggest extending the technique in this subsection to account quantitatively for multiple attributes. This discussion, however, will consider only the male/female distinction.

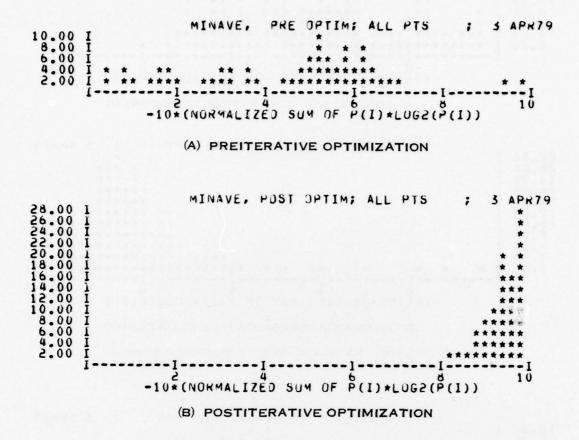
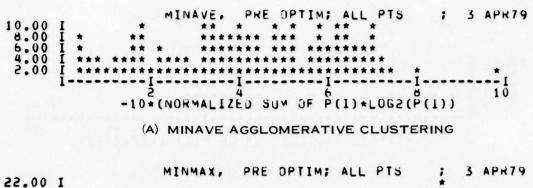


Figure 25. Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes for MINAVE Agglomerative Clustering of Recognition Patterns

The proposal is that the average information gained by knowing in which class a point falls should be reduced by the *a priori* knowledge of an attribute of that point, if the classes represent that attribute. Alternately stated, the proposal is that the average uncertainty about the class in which a point falls should be reduced by the amount of certainty gained about the class membership, knowing the attribute (the sex in this case). From the information theory literature (e.g., Reza³³), the average uncertainty is the entropy, given by*

$$H(c) \approx -\sum_{i=1}^{c} p(i) \log p(i)$$
 (21)

³³ F.M. Reza, An Introduction to Information Theory. New York: McGraw-Hill, 1961. *All logarithms are taken to base 2.



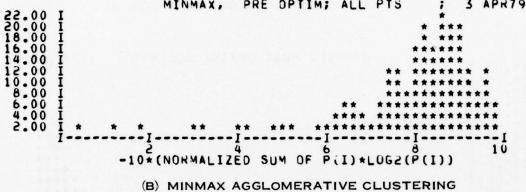
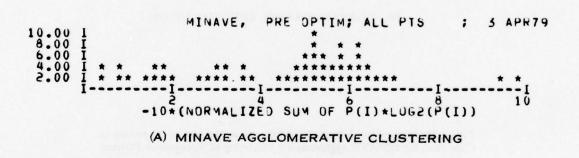
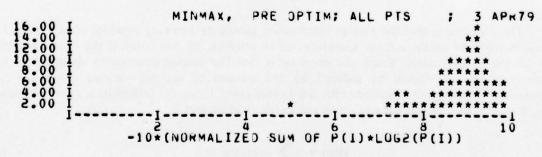


Figure 26. Histograms of Normalized Entropy as Measure of Dispersion in Class Size for Preiterative Optimization of Scanning Patterns





(B) MINMAX AGGLOMERATIVE CLUSTERING

Figure 27. Histograms of Normalized Entropy as Measure of Dispersion in Class Sizes for Preiterative Optimization of Recognition Patterns

TABLE 7. CONTINGENCY TABLE OF THE FIRST KIND* RESULTS FOR PREITERATIVE OPTIMIZATION

For Scanning Patterns

	Ref			Contin	gency Table	Ratio for 6	Given Numb	er of Classe	s	
Digit	Point	2	3	4	5	6	7	8	9	10
0	1	0.440	0.548	0.542	0.536	0.470	0.470	0.392	0.416	0.434
0	2	0.663	0.536	0.548	0.373	0.518	0.536	0.524	0.542	0.524
0	3	1.000	0.928	0.934	0.825	0.536	0.500	0.488	0.506	0.416
1	1	0.815	0.643	0.655	0.756	0.750	0.720	0.744	0.643	0.643
1	2	0.964	0.833	0.661	0.661	0.661	0.524	0.554	0.357	0.369
2	1	0.929	0.911	0.702	0.708	0.452	0.470	0.554	0.399	0.399
2	2	0.765	0.756	0.524	0.286	0.500	0.482	0.482	0.583	0.476
2	3	0.613	0.821	0.821	0.714	0.744	0.738	0.732	0.804	0.798
3	1	0.858	0.657	0.639	0.509	0.633	0.609	0.533	0.568	0.586
3	2	0.592	0.432	0.438	0.456	0.456	0.604	0.639	0.544	0.533
4	1	0.916	0.940	0.725	0.689	0.533	0.533	0.401	0.467	0.467
4	2	1.000	0.641	0.641	0.641	0.461	0.461	0.551	0.617	0.605
5	1	0.615	0.538	0.544	0.527	0.509	0.669	0.657	0.663	0.675
5	2	0.609	0.396	0.391	0.391	0.556	0.604	0.556	0.651	0.639
6	1	0.970	0.707	0.719	0.551	0.527	0.587	0.611	0.557	0.563
6	2	0.569	0.575	0.389	0.587	0.593	0.581	0.479	0.473	0.533
6	3	0.665	0.563	0.671	0.599	0.539	0.581	0.539	0.551	0.599
7	1	0.994	0.756	0.583	0.417	0.429	0.429	0.417	0.393	0.399
7	2	0.881	0.821	0.607	0.500	0.518	0.435	0.536	0.548	0.548
7	3	0.827	0.536	0.524	0.571	0.494	0.494	0.589	0.565	0.530
8	1	0.798	0.512	0.464	0.494	0.589	0.423	0.393	0.446	0.440
8	2	0.583	0.679	0.571	0.601	0.565	0.589	0.417	0.417	0.458
9	1	0.515	0.544	0.562	0.373	0.527	0.462	0.485	0.497	0.544
9	2	0.621	0.538	0.538	0.580	0.592	0.621	0.627	0.633	0.716

For Recognition Patterns

	Contingency Table Ratio for Given Number of Classes												
Digit	2	3	4	5	6	7	8	9	10				
0	0.807	0.578	0.590	0.434	0.434	0.416	0.440	0.434	0.349				
1	0.750	0.607	0.619	0.530	0.435	0.363	0.387	0.375	0.369				
2	0.500	0.500	0.482	0.494	0.542	0.536	0.512	0.524	0.411				
3	0.633	0.556	0.408	0.604	0.544	0.491	0.509	0.456	0.479				
4	0.772	0.790	0.790	0.515	0.491	0.491	0.497	0.491	0.497				
5	0.858	0.775	0.763	0.728	0.604	0.633	0.633	0.491	0.491				
6	0.778	0.707	0.605	0.617	0.593	0.587	0.599	0.587	0.689				
7	0.494	0.821	0.750	0.565	0.583	0.583	0.595	0.560	0.583				
8	0.554	0.464	0.631	0.619	0.518	0.524	0.548	0.464	0.548				
9	0.491	0.538	0.556	0.574	0.396	0.361	0.320	0.408	0.432				

^{*}Refer to Subsection IV.D.3.

TABLE 8. CONTINGENCY TABLE OF THE FIRST KIND* RESULTS FOR POSTITERATIVE OPTIMIZATION

For Scanning Patterns

	Ref			Contin	gency Table	e Ratio for	Given Numb	er of Classe	S	
Digit	Point	2	3	4	5	6	7	8	9	10
0	1	0.633	0.964	0.657	0.795	0.699	0.627	0.633	0.651	0.735
0	2	1.000	1.000	0.723	0.807	0.590	0.596	0.596	0.741	0.753
0	3	1.000	0.669	1.000	0.614	0.687	0.699	0.633	0.669	0.657
1	1	1.000	0.857	1.000	0.863	0.631	0.458	0.518	0.732	0.827
1	2	0.958	0.815	0.554	0.702	0.548	0.708	0.655	0.708	0.702
2	1	0.917	0.815	0.625	0.708	0.708	0.685	0.655	0.571	0.571
2	2	1.000	1.000	0.637	0.726	0.857	0.810	0.738	0.732	0.714
2	3	1.000	0.696	1.000	0.571	0.756	0.917	0.869	0.935	0.899
3	1	0.675	0.775	0.462	0.799	0.917	0.663	0.568	0.663	0.669
3	2	0.988	0.633	0.970	0.716	0.686	0.716	0.734	0.710	0.704
4	1	1.000	1.000	0.856	0.826	0.587	0.689	0.695	0.766	0.641
4	2	1.000	0.808	0.665	0.665	0.605	0.641	0.719	0.743	0.635
5	1	0.675	0.503	0.923	0.799	0.893	0.716	0.462	0.793	0.586
5	2	1.000	0.964	0.959	0.846	0.710	0.941	0.822	0.675	0.675
6	1	1.000	1.000	0.844	0.707	0.880	0.850	0.665	0.599	0.635
6	2	1.000	0.647	0.850	0.737	0.814	0.665	0.575	0.713	0.713
6	3	1.000	0.467	0.689	1.000	0.886	0.629	0.701	1.000	0.904
7	1	1.000	0.988	0.500	0.929	0.750	0.601	0.720	0.655	0.726
7	2	1.000	0.577	0.762	0.875	0.964	0.821	0.702	0.607	0.667
7	3	0.488	0.988	0.982	0:601	0.613	0.714	0.661	0.798	0.756
8	1	1.000	0.839	0.940	0.685	0.488	0.470	0.607	0.690	0.661
8	2	0.952	0.685	0.637	0.673	0.679	0.571	0.851	0.714	0.690
9	1	1.000	0.781	0.680	0.645	0.787	0.775	0.686	0.746	0.917
9	2	0.692	0.757	0.846	1.000	0.746	0.716	0.799	0.805	0.870

For Recognition Patterns

Contingency Table Ratio for Given Number of Classes

Digit	2	3	4	5	6	7	8	9	10
0	1.000	0.753	0.789	0.590	0.639	0.861	0.783	0.753	0.614
1	0.940	0.738	0.940	0.780	0.655	0.661	0.792	0.827	0.685
2	1.000	0.964	0.798	0.851	0.708	0.744	0.696	0.565	0.607
3	1.000	0.645	0.686	0.562	0.556	0.515	0.562	0.562	0.592
4	0.988	0.916	0.491	0.617	0.563	0.743	0.587	0.563	0.647
5	1.000	0.834	0.728	0.817	0.609	0.580	0.734	0.757	0.852
6	1.000	0.982	0.743	0.796	0.832	0.689	0.796	0.629	0.593
7	1.000	0.929	0.726	0.530	0.679	0.554	0.714	0.655	0.601
8	1.000	1.000	1.000	0.958	0.839	0.679	0.685	0.565	0.726
9	1.000	0.793	1.000	0.935	0.609	0.669	0.645	0.651	0.669

^{*}Refer to Subsection IV.D.3.

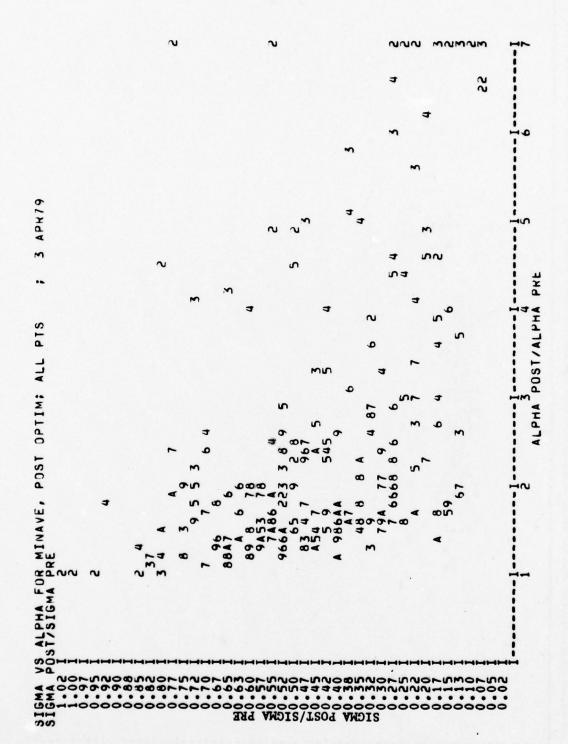
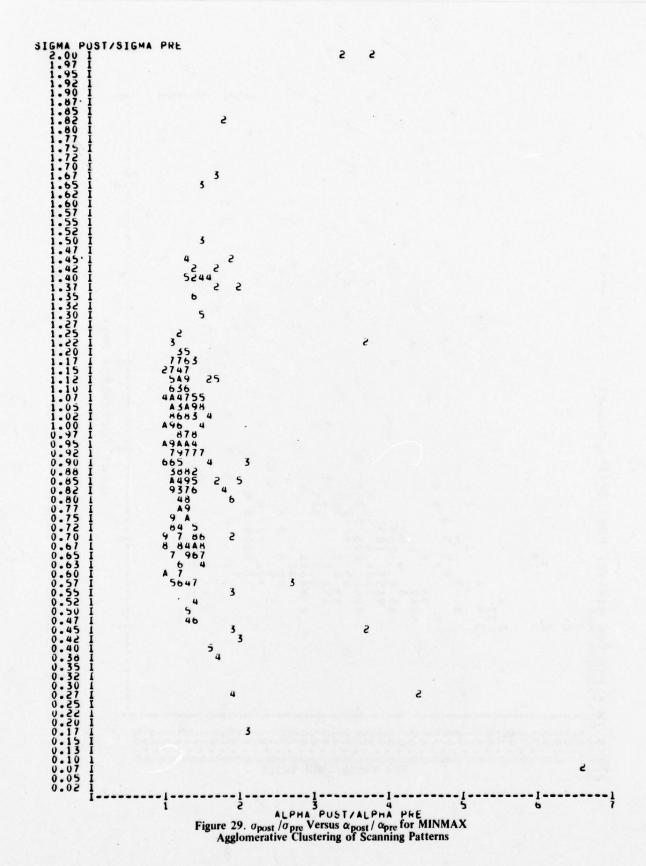


Figure 28. σ_{post}/σ_{pre} Versus α_{post}/α_{pre} for MINAVE Agglomerative Clustering of Scanning Patterns



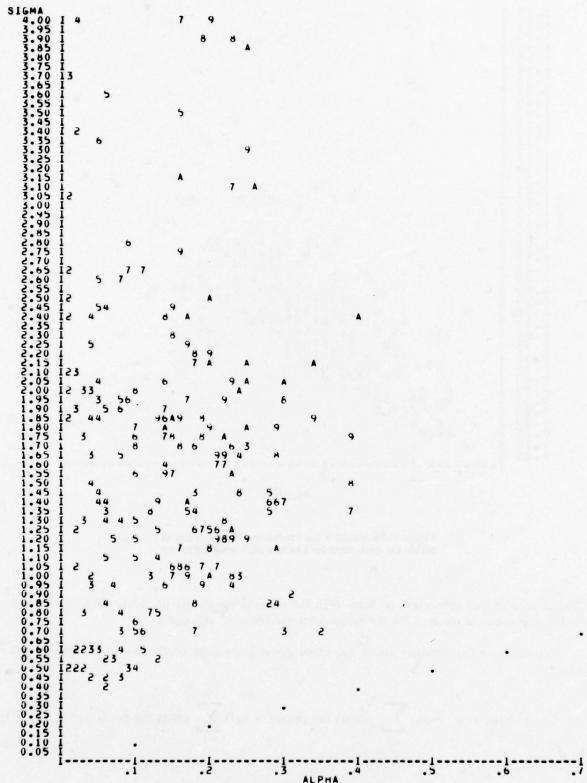


Figure 30. σ Versus α for Preiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns

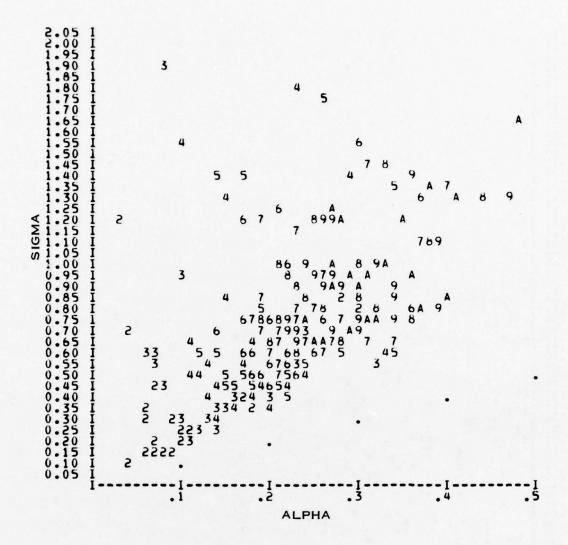


Figure 31. σ Versus α for Preiterative Optimization of MINMAX Agglomerative Clusters of Scanning Patterns

("H" is used in this subsection to agree with the information theory literature. The "E" used in the last subsection is reserved for the normalized entropy, i.e., H/log₂ c.)

The average information about the class, given knowledge of the sex, is the conditional entropy:

$$H(c|s) = -p(m) \sum_{i=1}^{c} p(i|m) \log p(i|m) - p(f) \sum_{i=1}^{c} p(i|f) \log p(i|f)$$
 (22)

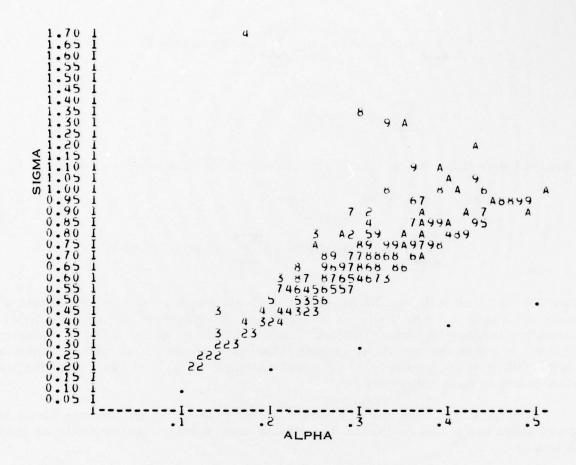


Figure 32. σ Versus α for Postiterative Optimization of MINAVE Agglomerative Clusters of Scanning Patterns

The average information gained from the clustering is then reduced by the *a priori* information represented in the classes, yielding what information theorists term I (c:s), the average of the mutual information between the class and the sex. The value for I(c:s) is given by

$$I(c;s) = H(c) - H(c|s)$$

$$= -\sum_{i=1}^{c} p(i) \log p(i) + p(m) \sum_{i=1}^{c} p(i|m) \log p(i|m)$$
 (23)

+
$$p(f)$$
 $\sum_{i=1}^{c} p(i|f) \log p(i|f)$

Note that since H(c) - H(c|s) = H(s) - H(s|c), the term I(c;s) can also be written as:

$$I(c;s) = -p(m) \log p(m) - p(f) \log p(f)$$

$$+ \sum_{i=1}^{c} p(i) p(m|i) \log p(m|i) + \sum_{i=1}^{c} p(i) p(f|i) \log p(m|i)$$
(24)

Note also that I(c;s) = H(c) for the two-class case when a population with an equal number of males and females are divided into an all-male class and an all-female class since the class is uniquely determined by knowing the sex. This corresponds to what is called a "noise-free channel" in information theory. In contrast, when no information is transmitted through a channel, H(c) = H(c|s), yielding I(c;s) = 0, which corresponds to the case of each class having an equal number of males and females.

However, I(c;s) is used here as a measure of the "goodness" of the clustering relative to the condition (sex in this case) tested. The estimates used for the various probabilities are given in terms of

n	Total number of samples		
n _i	Number of samples in class i		
$n^{\mathbf{m}}$	Number of males		
n^f	Number of females		
n _i ^m	Number of males in class i		
n!	Number of females in class i		

The given probabilities are estimated as

$$p(i) = n_i/n$$

$$p(m) = n^m/n$$

$$p(f) = n^f/n$$

$$p(i|m) = n_i^m/n^m$$

$$p(i|f) = n_i^f/n^f$$

$$p(m|i) = n_i^m/n_i$$
$$p(f|i) = n_i^f/n_i$$

Hence, I(c;s) is calculated by

$$I(c;s) = -\sum_{i=1}^{c} (n_i/n) \log (n_i/n)$$

$$+(n^m/n) \sum_{i=1}^{c} (n_i^m/n^m) \log (n_i^m/n^m)$$

$$+(n^f/n) \sum_{i=1}^{c} (n_i^f/n^f) \log (n_i^f/n^f)$$
(25)

Values of 1 are given in Tables 9 through 12 for scanning patterns and in Tables 13 through 16 for recognition patterns, along with another measure, R, the distribution of males and females among the classes, for the reader with a less esoteric inclination. The value for R is given by

$$R = \frac{1}{n} \sum_{i=1}^{c} \min (n_i^f, n_i^m)$$
 (26)

which is a measure of the residue in each of the classes. The tables in each of the two sets are arranged by clustering algorithm in the following order:

- (1) MINAVE agglomerative clustering; preiterative optimization
- (2) MINAVE agglomerative clustering; postiterative optimization
- (3) MINMAX agglomerative clustering; preiterative optimization
- (4) MINMAX agglomerative clustering; postiterative optimization.

The information in Tables 9 through 12 is summarized in Figure 33 with histograms of I for each of cases 1 through 4 above. Likewise, the information in Tables 13 through 16 is summarized in Figure 34. It is clear from these two figures that iterative optimization to minimize J_e improves the resulting clusters, assuming the male/female distinction is valid (and from an acoustic-phonetic standpoint, it is). In addition, these two figures show that MINMAX agglomerative clustering yields better clusters before the iterative optimization than does MINAVE. However, no clear preference results between the MINAVE and MINMAX clusters after iterative optimization, showing that the iterative optimization algorithm was robust enough to produce good clusters from either the MINAVE or MINMAX agglomerative clustering, even though the starting partition produced from the MINAVE agglomeration was clearly inferior.

TABLE 9. MUTUAL INFORMATION AND RESIDUES FOR PREITERATIVE OPTIMIZATION OF MINAVE AGGLOMERATIVE CLUSTERS OF SCANNING PATTERNS

01	3/1	RP (c: 2	3	4	5	6	7	8	9	10
0	1	I:	0.030	0.091	0.119	0.179 0.355	0.195	0.217	0.288	0.289	0.289
0	Ş	I:	0.006	0.024	0.036	0.038	0.372	0.372	0.395 0.175	0.401	0.466
0	3	I:	0.052	0.093	0.103	0.103	0.125	0.125	0.130	0.143	0.143
1	1	I:	0.006	0-014		0.034		0.049		0.075	0.075
1	5	I:	0.006	0.012			0.086	0.086	0.097	0.097	0.124
5	1	1:	0.005			0.032		0.039		0.064	0.071
3	ş	I:	0.018	0.018		0:025	0.714		0.755	0.757	0.794
5	3	I:	0.007		0.026	0.030				0.069	
3	1	I:	0.000	0.028	0.035				0.054	0.079	0.081
3	Ş	I:	0.000	0.007	0.046	0.046	0.074	0.154	0.154	0.162	0.165
4	1	I:	0.006	0.049			0.118			0.133	0.138
4	5	I: R:	0.006	0.018	0.018	0.124	0.140	0.146	0.208	0.265	0.320
5	1	I:	0.055	0.146	0.218 0.302	0.222	0.230	0.236	0.236	0.245	0.245
5	Ş	I:	0.000	0.006	0.043	0.049	105.0	0.314	0.332	0.356	0.358
6	1	I:	0.206	0.230	0.230	0.268	0.272	0.620	0.621 0.078	0.642	0.659
6	Ş	I: R:	0.006	0.012	0.018	0.633	0.634	0.634	0.714	0:714	0.720
6	3	I:	0.005	0.017	0.051	0.178	0.185	0.207	0.207	0.248	0.248
77	1	I:	0.023	0.029	0.570	0.598	0.620	0.623	0.623	0.641	0.648
7	5	I:	0.006	0.591	0.596	0.601	0.603	0.604	0.604	0.604	0.604
7	3	I: R:		0.031							
8	1	I: R:	0.041	0:104	0:111	0:139	0.405	0.421	0.475	0.476	0.478
8	5	I:	0.008	0.082	0.103	0:113	0:126	0.126	0.126	0:126	0.144
9	i	I:	0.023	0.046	0.050	0.050	0.270	0.285	0.306	0.334	0.353
9	Ş	I:	0:012	0.021	0.032	0.034	0.040	0.082	0.085	0.099	0.247

TABLE 10. MUTUAL INFORMATION AND RESIDUES FOR POSTITERATIVE OPTIMIZATION OF MINAVE AGGLOMERATIVE CLUSTERS OF SCANNING PATTERNS

DG/RP C: 2	3	4	5	6	7	8	9	10
0 1 I: 0.0 0 1 R: 0.4	34 0.575	0.387	0.524	0.588	0.552	0.598	0.621	0.631
0 2 I: 0.3 0 2 R: 0.1	89 0.447 57 0.181	0.404	0.477	0.493	0.560	0.556 0.145	0.003	0.578 0.139
0 3 I: 0.0 0 3 R: 0.3	9 0.067	0.177	0.245	0.286	0.215	0.257	0.269	0.276
1 1 I: 0:0 1 1 R: 0:4	23 0.015	0.040	0.048	0.059	0.046	0.051	0.079	0.102
1 2 I: 0.3 1 2 R: 0.1	23 0.285	0.455	0.428	0.420	0.382	0.464	0.401	0.416
2 1 I: 0.0 2 1 R: 0.3	75 0.063 57 0.375	0.089	0.080 0.363	0.129	0.126	0.132	0.134	0.188
5 5 I: 0.6	74 0.658	0.569	0.705	0.712	0.759 0.083	0.732	0.665	0.704
2 3 I: 0.0 2 3 R: 0.4	3 0.013 5 0.435	0.014	0.069	0.056 0.375	0.069	0.065	0.113 0.345	0.155 0.333
3 1 I: 0.0 3 1 R: 0.4	0.029	0.045	0.079 0.355	0.057 0.373	0:131 0:337	0.180 0.308	0.169	0.158
3 2 I: 0.0 3 2 R: 0.3	58 0.062 51 0.361	0.103	0.159	0.141	0.157	$0.181 \\ 0.314$	0.281	0.289
4 1 I: 0:0 4 1 R: 0:4	0.036	0.306	0.296	0.309	0.286	0.258	0.402	0.355
4 2 I: 0:2 4 2 R: 0:2	0 0.341 34 0.198	0.343	0.256	0.316 0.234	0.342	0.395	0.363	0.435
5 1 I: 0:5 5 1 R: 0:1	23 0.435	0.397	0.522	0.610	0.617 0.130	0.451	0.557	0.535
5 2 I: 0:4 5 2 R: 0:1	35 0.361 36 0.172	0.416	0.329	0.429	0.378	0.424	0.400	0.428
6 1 I: 0:2	55 0.451 9 0.162	0.426	0.614	0.600	0.636	0.723	0.671	0.722
6 2 I: 0.7 6 2 R: 0.0	19 0.761	0.799	0.740	0.668	0.633	0:732	0.603	0.729
6 3 I: 0.1 6 3 R: 0.2	77 0.236 59 0.257	0.259	0.346	0.304	0.413	0.413	0.345	0.343
7 1 1: 0:7 7 1 R: 0:0	0.599 8 0.089	0.701	0.683	0.728	0.655	0.662	0.650	0.764
7 2 I: 0:6 7 2 R: 0:0	0.685 77 0.071	0.573	0.720	0.680	0:717	0:709	0:736	0.660
7 3 I: 0.0 7 3 R: 0.4	4 0.335 5 0.238	0.346	0.313	0.375	0.425	0.404 625.0	0.388	0.445
8 1 I: 0:1 8 1 R: 0:2	39 0.281 6 0.214	0.183	0.300	805.0	0.322	0.461	0:417	0.427
8 2 I: 0:3 8 2 R: 0:1	34 0.238 5 0.250	0.289	0:575	392.0	0.528	0.250	0.496	0.527
3 1 A: 0:4	2 0:515	0.484	0:456	0.544	0.514	0.557	0.566	0.568
9 2 I: 0:0 9 2 R: 0:4	5 0:121	0:296	0:292	0:252	0.263	0.281	0:291	0:333

TABLE 11. MUTUAL INFORMATION AND RESIDUES FOR PREITERATIVE OPTIMIZATION OF MINMAX AGGLOMERATIVE CLUSTERS OF SCANNING PATTERNS

DG/RP C: 2	3 4	5	6	7	8	9	10
0 1 1: 0.222 0 1 R: 0.235	0.356 0.398	0.398	0.399	0.423	0.423	0.438	0.440
0 2 I: 0.022 0 2 R: 0.416	0.105 0.133 0.380 0.380	0.373	0.390	0.453	0.461	0.461	0.466
0 3 I: 0.052 0 3 R: 0.392	0.060 0.072	0.079	0:165	0.182	0.277	0.291	0.291
1 1 I: 0.009 1 1 R: 0.464	0.017 0.032	0.035	0.092	0.119	0.120	0.130	0.134
1 2 I: 0.031 1 2 R: 0.458	0.032 0.155 0.458 0.310	0.158 0.310	0.159	0.214	0.223	0.525	0.525
2 1 I: 0.005 2 1 R: 0.476	0.006 0.065 0.476 0.363	0.080 0.357	0.088	0.104	0:104 0:357	0.104	0.126
2 2 I: 0.005 2 2 R: 0.464	0.005 0.343	0.455	0.493	0.516	0.522	0.528	0.528
2 3 I: 0.002 2 3 R: 0.462	0.006 0.033 0.458 0.429	0.034	0.041	0.054	0.056 0.411	0.057 0.411	0.063
3 1 I: 0.004 3 1 R: 0.467	0.078 0.090 0.367 0.349	0.173 0.306	0.174	0.203	0.204	0.204	0.204
3 2 I: 0.093 3 2 R: 0.331	0.133 0.169	0.169	0.199	0.207 272	0:212 272	0.221	0.230
4 1 I: 0.000 4 1 R: 0.497	0.026 0.186 0.473 0.293	0.239	0:311	0.350	0.361 0.198	0.363	0.370 0.198
4 2 I: 0.006 4 2 R: 0.497	0.013 0.062 0.449 0.395	0.084	0.099	0.130 0.359	0.173 0.329	0.405	0.420
5 1 I: 0.001 5 1 R: 0.479	0.083 0.144 0.396 0.361	0.167	0.304	0.340	0.382	0.385	0.418
5 2 I: 0.065 5 2 R: 0.355	0.266 0.290	0.375	0.387 0.189	0.414	0.470	0.483	0.483
6 1 I: 0.250 6 1 R: 0.251	0.330 0.352 0.246	0.556	0.558	0.569	0.569	0.570	0.592
6 2 I: 0.550 6 2 R: 0.102	0.556 0.619	0.621	0.622	0.691	0.697	0:712	0.715
6 3 I: 0:111 6 3 R: 0:323	0.169 0.201	0.212	0.230	0.234	0.240	0.247	0.265
7 1 I: 0.017 7 1 R: 0.488	0.185 0.213	0.260	0.285	0.355	0.323	0.444	0.455
7 2 I: 0.008 7 2 R: 0.458	0.523 0.524	0.534	0.582	0.610	0.631	0.655	0.655
7 3 I: 0.198 7 3 R: 0.321	0.246 0.251	0.251	565.0	995.0	0.266	0.279	0.295
8 1 I: 0.081 8 1 R: 0.351	0.129 0.136	0.142	0.143	0.191	0.261	0.267	785.0
8 2 I: 0.095 8 2 R: 0.327	0.200 0.217	0.224	0.224	0.231	0.487	0.496	0.505
9 1 I: 0.220 9 1 R: 0.231	0.236 0.271	0.380	0.396	0.523	0.548	0.548	0.553
9 2 I: 0.059 9 2 R: 0.361	0.063 0.080	0:113	0.285	0.265	0.265	0.294	0.328

TABLE 12. MUTUAL INFORMATION AND RESIDUES FOR POSTITERATIVE OPTIMIZATION OF MINMAX AGGLOMERATIVE CLUSTERS OF SCANNING PATTERNS

DG/RP C: 2	3 4	5	6	7	8	9	10
0 1 I: 0.574 0 1 R: 0.090	0.666 0.609	0.631	0.630	0.627	0.654	0.659	0.665
0 2 I: 0.389 0 2 R: 0.157	0.447 0.411	0.492	0.534	0.542	0.536	0.604	0.645
0 3 1: 0.049 0 3 R: 0.392	0.133 0.177	0.147	0.191	0.217	0.287 0.235	0.353	0.403
1 1 I: 0.023 1 1 R: 0.411	0.035 0.040	0.065	0.087 0.357	0.156 0.315	0.135 0.345	0.130	0.137
1 2 I: 0.328 1 2 R: 0.196	0.290 0.459	0.443	0.367	898.0	0.339	0.452	0.424
2 1 I: 0.053 2 1 R: 0.393	0.077 0.062	0.096	0:061 0:387	0.066	0.286	0.148	0.176
2 2 I: 0.674 2 2 R: 0.060	0.658 0.711	0.617	0.684	0.792 0.077	0.698	0.709	0.751
2 3 I: 0.013 2 3 R: 0.435	0.019 0.014	0.021	0.077 0.387	0.093	0.096	0:128	0.145
3 1 I: 0.002 3 1 R: 0.473	0.012 0.023	0.076 0.379	0.066	0.100 0.355	0.082 0.373	0.074	0.143
3 2 I: 0.049 3 2 R: 0.373	0.108 0.093 0.325 0.349	0.095	0.095	0.158	0.175	0:179	0.214
4 1 I: 0.001 4 1 R: 0.485	0.036 0.287	0.305	0.265	0:255	0.298	0.294	0.316
4 2 I: 0.220 4 2 R: 0.234	0.212 0.263	0.286	0.315	0.347	855.0	0.367	0.376
5 1 I: 0:044 5 1 R: 0:379	0.414 0.457 0.176 0.148	0.540	0.627	0.552	0.546	0.589	0.554
5 2 I: 0.435 5 2 R: 0.136	0.408 0.393 0.166 0.183	0.380	0.439	0.364	0.397	0.532	0.536
6 1 I: 0:235 6 1 R: 0:269	0.451 0.619	0.625	0.612	0:611	0.600	0.558	0.607
6 2 I: 0.749 6 2 R: 0.042	0.754 0.762				0.705	0.760	0.790
6 3 I: 0.177 6 3 R: 0.269	0.194 0.285	0.346	0.390	0.381	0.314	0.344	0.368
	0.601 0.624				0.145	0.781	
7 2 1: 0:607 7 2 R: 0:077	0.487 0.661 0.131 0.077	0.665	0.684	0.684	0.747	0.723	0.723
	0.332 0.347						
	0.185 0.236						
	0.452 0.366 0.190 0.244						
	0.449 0.534 0.130 0.136						
9 2 1: 0:110 9 2 R: 0:320	0.030 0.240	0.292	0.288	0.270	0.291	0.302	0.367

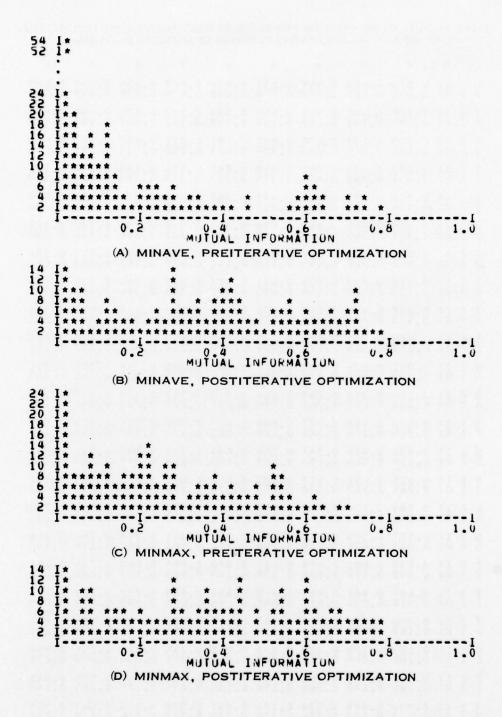


Figure 33. Comparison of Mutual Information for Clustered Scanning Patterns Using Four Clustering Algorithms

TABLE 13. MUTUAL INFORMATION AND RESIDUES FOR PREITERATIVE OPTIMIZATION OF MINAVE AGGLOMERATIVE CLUSTERS OF RECOGNITION PATTERNS

DG/	RP C	:: 5	3	4	5	6	7	8	9	10
										0.693
								0.497		
								0.729		
								0.572		
4 0	I:	0.006	0.018 0.485	0.024	0:197 0:311	0.198	0:196	0.235	0.251	0.278
								0.768		
								0.640		
								0.793		
8 0	I:	0.006	0.023	0.872 0.018	0.873 0.018	0.873 0.016	0.874	0.874 0.018	0.875	0.881
9 0	I:	0.031	0.119	0.128	0.128 0.385	0.136 0.367	0.170	0:170	0.391	0.393

TABLE 14. MUTUAL INFORMATION AND RESIDUES FOR POSTITERATIVE OPTIMIZATION OF MINAVE AGGLOMERATIVE CLUSTERS OF RECOGNITION PATTERNS

0	· E · · · · · · · · · · · · · · · · · ·			
DG/RP C: 2	3 4	5 6	7 8	9 10
0 0 I: 0.78	0.678 0.685 0.060 0.060	0.688 0.656	0.810 0.743 0.042 0.054	0.696 0.788
1 0 I: 0:34 1 0 R: 0:17	11 0.318 0.434	0.447 0.446 0.179 0.155	0.375 0.560 0.179 0.107	0.501 0.490 0.143 0.149
2 0 1: 0.66 2 0 R: 0.06	26 0.515 0.670 33 0.131 0.083	0.676 0.659	0.631 0.768 0.125 0.054	0.717 0.800
3 0 I: 0.40 3 0 R: 0.19	01 0.347 0.404 0.219 0.148	0.483 0.421 0.160	0.561 0.612 0.124 0.130	0.623 0.627
4 0 I: 0.15 4 0 R: 0.28	51 0:165 0:213 51 0:311 0:311	0.298 0.398 0.234 0.180	0.396 0.322	0.325 0.432
5 0 I: 0.54 5 0 R: 0.09	19 0.583 0.829 5 0.142 0.063	0.795 0.706 0.089 0.118	0.724 0.725	0.700 0.711
6 0 I: 0.84 6 0 R: 0.08	11 0.769 0.682	0.728 0.734 0.060	0.739 0.785	0.811 0.813
7 0 I: 0.95 7 0 R: 0.00	3 0.829 0.819 6 0.030 0.077	0.768 0.848 0.065 0.036	0.852 0.933	0.785 0.791
8 0 I: 0.70 8 0 R: 0.05	3 0.797 0.840 4 0.036 0.030	0.840 0.832 0.030 0.042	0.764 0.795	0.802 0.887
	54 0.642 0.505 5 0.077 0.154			

TABLE	15. MUTUAL	INFORMATION AND RESIDUES FOR PREITERATIVE OPTIMIZATION
	OF MINMAX	AGGLOMERATIVE CLUSTERS OF RECOGNITION PATTERNS

DG	/R	P	: 2	3	4	5	6	7	8	9	10
0	0	I: R:	0.212 106.0	0.707 0.054	0.711 0.054	0.716	0.721	0.721	0.724 0.054	0.740	0.776
									0.433 0.179		
									0.549 0.149		
3	0	I: R:	0.325	0.326	0.356	0.407	0.428	0.473	0.510 0.130	0.538 0.130	0.538 0.130
4	0	I: R:	0:162	0:170 0:311	0.204	0.205	0.205	0.230	0.250	0.250	0.250
									0.797 0.041		
6	0	I: R:	0.204	0.655	0.674	0.686	0.692	0.692	0.692	0.713	0.713
77	0	I: R:	0.953	0.963	0.963	0.963	0.969	0.969	0.969	0.969	0.969
8	0	I:	0.558	0.606	0.707	0.707	0.707	0.730	0.736 0.089	0.736	0.738
9	0	I:	0.369	0.451	0.453	0.487	0.502	0.502	0.574 0.130	0.582	0.595

TABLE 16. MUTUAL INFORMATION AND RESIDUES FOR POSTITERATIVE OPTIMIZATION OF MINMAX AGGLOMERATIVE CLUSTERS OF RECOGNITION PATTERNS

D	G/	RP	C: 5	3	4	5	6	7	8	9	10
0	0	I: R:	0.787	0.689	0.678	0.726	0.720	0.764	0.709	0.719	0.734
1	0	I:	0.310	0:409	0.419	0.362	0.396	0.523	0.596	0.515	0.516
5	0	I:	0.628	0.466	0.611	0.710	0.690	0.664	0.675	0.795	0.789
3	0	I:	0.401	0.443	0.473	0.579	0.602	0.602	0.557	0.587	0.620
4	0	1: R:	0.153	0.176	0.252	0.239	0.294	0.325	0.370	0.448	0.416
5	0	I:	0.549	0.735	0.670	0.693	0.757	0.742	0.804	0.805	0.793
6	0	I:	0.841	0.770	0.764	0.745	0.730	0.736	0.616	0.829	0.812
7	0	I: R:	0.953	0.911	0.828	0.887	0.888	0.881	0.669	0.936	0.935
8	0	I:	0.703	0.797	0.840	0.840	0.800	0.823	0.634	0.828	0.878
9	00	I:	0.554	0.494	0.505	0.700	0.568	0.519	0.645	0.631	0.704

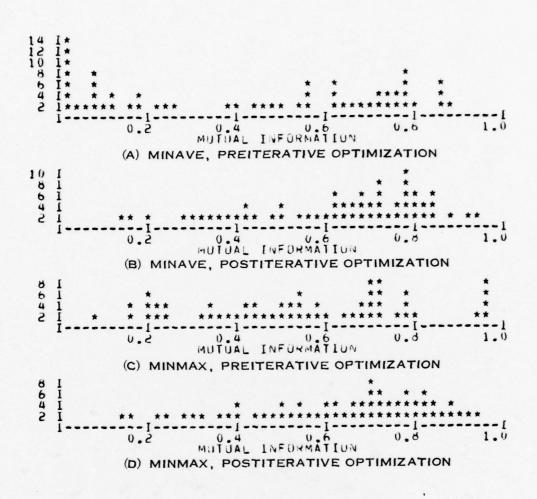


Figure 34. Comparison of Mutual Information for Clustered Recognition Patterns Using Four Clustering Algorithms

SECTION V GENERAL-PURPOSE SPEECH I/O CAPABILITY

This section describes the methods used in the AP-120B version of the digit recognition program. Topics covered include the data collection hardware, the filter simulation, the auto-correlation computation, and a discussion of digitizing and playback utilities.

A. SYSTEM DESCRIPTION

The primary impetus of designing this remote speech I/O facility was to relieve the host from the burden of controlling analog-to-digital (A/D) and digital-to-analog (D/A) converters. A second consideration was to develop a method to get data into an array processor at a high rate of speed, to allow real-time data collection and processing. Figure 35 is a block diagram of the resulting speech I/O subsystem.

The configuration consists of a TI 980B host computer, a Floating Point Systems AP-120B array processor, and a TI 990/10 computer with attached A/D and D/A converters. The 990/10 collects and plays out data under the control of a "mailbox" memory location in the AP-120B. Therefore, the 990/10 can be controlled by either the 980B or the AP-120B. The AP-120B is used primarily to reduce the quantity of data by transforming the raw speech to a more compact form (e.g., filtering or preprocessing). In a typical application, the host would request data from the 990/10, request the AP-120B to process the data, and then request that the results of that processing be sent to the 980B.

Software directly used in the I/O subsystem consisted of two parts. The first is the 990/10 software that controls the A/D and D/A, buffers the input and output speech, and controls speech I/O to the AP-120B. The second piece of software runs on the 980B and is basically a device driver for the 990/10. This driver handles the channel protocol as well as all I/O between the 980B and the AP-120B. In addition, existing software that digitizes and edits speech data using the 980B internal A/D and D/A was modified to use the new data acquisition subsystem.

The offloading of the host was carried one step further in the digit-recognition program. In order to free the host from controlling the 990/10 and AP-120B, the AP-120B was put in charge of this entire process. Since the AP-120B has direct memory access (DMA) capability to the

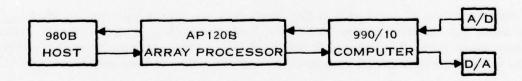


Figure 35. Speech Channel Block Diagram

host, the host need only tell the AP-120B where to place the processed data and when to begin collection. This scheme allows continuous input of speech, since host activity is totally disjointed from the data collection or processing. With all this computational burden removed from the host, it can easily keep pace with real-time processing.

B. FILTER SIMULATION IN THE AP-120B

The first process initiated on the input speech data is a recursive filter simulation. The center frequencies and bandwidths of this simulation are designed to match those of the hardware filters specified in Subsection II.A and Appendix A. The filter model included both preemphasis and envelope shaping to match the hardware filters. These 16 filters were typically sampled every 10 milliseconds. This filter simulation accounts for approximately 60 percent of real time when data are collected at an 80-microsecond sample rate.

The output of the filter simulator is then preprocessed, and the output of the preprocessor is sent to the 980B host memory. The preprocessing is the same as that described in Appendix A.

C. VOICING DECISION FROM THE AP-120B AUTOCORRELATION PITCH TRACKER

An estimate of voicing was included in this speech input subsystem by performing an autocorrelation on the input speech. This consisted of sliding a window of speech over previous speech. Given a frame of speech data consisting of N samples, the last M samples of the frame, W_r , are used as a sliding window for a reverse correlation. The normalized cross-correlation, $R_r(K)$, between the sliding window, W_r , and the M speech samples earlier in time starting at the $(M - K_{min})th$ sample is used for this reverse correlation; i.e.,

$$R_{r}(K) = \frac{\left[\sum_{m=1}^{M} X (M - m + 1) X (M - m - K + 1)\right]}{\left(\left[\sum_{m=1}^{M} X^{2} (M - m + 1)\right]\left[\sum_{m=1}^{M} X^{2} (M - m - K + 1)\right]\right)}$$

$$(K_{min} \le K \le K_{max})$$
(28)

The maximum value of the $R_r(K)s$ is then selected as the voicing indicator. An $|R_r(K)_{max}|$ less than about 0.6 indicates an unvoiced frame, while an $|R_r(K)_{max}|$ approaching 1 indicates a strongly voiced frame. The value of K for the $R_r(K)_{max}$ was not used, although it corresponds to the value of the pitch period in samples. Typical values used in these calculations are:

M = 375 (375 samples at 80 microseconds = 30 milliseconds)

N = 125 (125 samples at 80 microseconds = 10 milliseconds)

 $K_{min} = 25$ (25 samples at 80 microseconds = 2 milliseconds)

 K_{max} = 250 (250 samples at 80 microseconds = 20 milliseconds)

SECTION VI EXPERIMENTAL RESULTS

A. SPEAKER-INDEPENDENT DIGIT RECOGNITION

Data Sets

A wide variety of digit-recognition testing was done using two types of data sets. The most heavily used data set was the total-voice evaluation data set. This test data set was part of a data base collected in a sound booth over a 3-month period at Texas Instruments. The test data set was extracted from around the middle of this 3-month period, to avoid the initial microphone fright of the subjects, and consisted of one repetition of one of 10 possible sets of 10 six-digit sequences uttered by 106 subjects (64 males, 42 females). The actual sequences used in this data collection are shown in Table 17. The test data were digitized, edited, and preprocessed during the total voice contract to ensure the precise replicability of the test data. However, since these "test" data were used for multiple experiments to evaluate the effect of parameter variations, the validity of the absolute recognition results is not assured. In addition, the data were also idealized by editing, which avoids spurious false recognition of background noises as true data. For these reasons, further experiments were performed on a second data base.

The second data base used is a subset of a large digit-recognition data base currently being collected in the speech community. The data being collected use sequences devised by Martin and Herscher, modified to include both the "oh" and the "zero" pronunciations of the digit zero. The texts being used in these data collections are shown in Table 18. All the multiple-digit sequences are supposed to be said in a continuous manner, although not all subjects always complied.

2. Digit-Recognition Results for Six-Digit Sequences

A total of 29 evaluation runs were made on the 1,060 six-digit sequences from the total voice evaluation data set. The overall digit recognition rates and conditions for all runs are given in Table 19. Even though no syntactic constraints (except length) were applied during the digit recognition, the total voice evaluation data set used was compatible with the design data since the following digit pairs, all nasal-to-vowel, glide or semivowel (or vice versa) transitions, were disallowed in both data sets:

0 - 1	1-8	2-9	3-9	4_9	9-1
0-8	2-1	3-1	4-1	7-1	9-8
0-9	2-8	3-8	4-8	7-8	

The first of these evaluation runs (no. 45), was the syntactically unconstrained digit recognition for the final evaluation for the total voice study. A detailed description of the thresholds and parameters for the evaluation runs is given in the total voice final report. The values of most of these parameters remained unchanged during this current study, except as noted in this section. These parameters are listed below along with the values used for run no. 45:

TABLE 17. THE 10 SETS OF 10 SIX-DIGIT SEQUENCES USED IN TESTING

061934 159034 253760 368405
253760
368405
430752
517943
675823
794302
852734
926034
026873
132057
234687
345768
403251
547602
651942
790234
879630
958170

Reference-point location parameters

Peak-to-valley ratio (PVR) = 1.10

Maximum valley point error (Max VPE) = 615

OPTSEQ (valley point sequencing parameters)

dt limits
$$(dt_{max}, dt_{min})$$
 see Table 11 of total voice final report¹ Expected dt $(d\hat{t})$

Minimum expected dt (used to determine dt* for the denominator in the point-pair error calculation:

$$dt^* = max (d\hat{t}, d\hat{t}_{min}) = 4$$

Time deviation weighting $(\beta) = 2$

Floor of valley point error (OFFSET) = 100

Hypothesized digit parameters

Minimum absolute average energy across recognition pattern (EN_{min}) = 150

Weighting of sequence error (SQ) contribution to total normalized error for digit $k(w_k) = 0.1$

TABLE 18. MARTIN-HERSCHER DIGIT TEXTS

Isolated Di	gits								
9	6	2 4	1 1	7	ZERO	3	5	8	ОН
9	6		1 1	7	ZERO	3	5	8	OH
Five-Digit	Codes								
(Pronounce	e "0" as	ZERO.)							
			08175	102	260	5580	6		
			67438	449	53	3214	6		
			29091	607	33	6863	0		
			91625	817	54	7924	1		
(Pronounce	e "0" as	OH.)							
			08175	102	260	5580	6		
			29091	607		6863	-		
There Diet	C-1								
Pronounce		ZERO.)							
(1)	525	(11)	990	(21)	531	(31)	005	(41)	033
(2)	759	(12)	583	(22) 3	349	(32)	140	(42)	477
(3)	101	(13)	171	(23) 5	565	(33)	819	(43)	680
(4)	626	(14)	098	(24) 1		(34)	974	(44)	306
	202	(15)		(25) 4		(35)	357	(45)	915
	727	(16)		(26) 8		(36)		(46)	
, ,	366	(17)		(27)		(37)		(47)	
	044	(18)		(28)		(38)		(48)	
, ,	843	(19)		(29)		(39)		(49)	
(10)		(20)		(30)		(40)		(50)	
(Pronounce	e "0" a	s OH.)							
		101	990	460) 0	05	033		
		202	098	076		40	680		
		044	670	370		08	306		

Normalizers to account for expected recognition error for digit k (TE_k) -values given later in this section.

Maximum allowable total normalized error for digit k (NE_k):

Digit 0 1 2 3 4 5 6 7 8 9

NE k 128 97 123 110 116 113 110 109 107 114

TABLE 19. SYNOPSIS OF EVALUATION RESULTS

Run No.	Percent Correct Recognition	Remarks (Changes From Previous Runs)
45	90.5	Baseline (TVBISS final parameters): Multiple reference patterns; minimum energy = 150; 0.1 percent NE thresholds; PVR = 1.1; maximum valley point error = 615; sequence length unconstrained.
46	90.0	Same as run no. 45 except new tree-searching subroutine (DECIDE) used with minimum separation = 3 centiseconds; maximum separation = 80 centiseconds.
47	89.4	Same as run no. 46 except sequence length constrained to 6.
48	89.9	Same as run no. 47 except maximum separation = 120 centiseconds and minimum energy = minimum [150, 0.1 (maximum energy of all hypothesized digits)].
49	90.0	Same as run no. 48 except PVR = 1.05 (for this run only).
50	90.3	Same as run no. 48 except point-pair error between reference points 1 and 3 added.
51	89.3	Same as run no. 50 except minimum energy = minimum [150, 0.2 (maximum energy of all hypothesized digits)].
52	91.8	Same as run no. 50 except total normalized error (NE) for 3 reference- point words multiplied by 0.9 (longer words are, in general, more reliably recognized).
53	91.9	Same as run no. 52 except minimum energy = minimum [150, 0.15 (maximum energy of all hypothesized digits)].
54	92.3	Same as run no. 53 except sequence length unconstrained (for this run only).
55	93.8	Same as run no. 53 except 3-bit quantized T-function added to scanning patterns and maximum valley-point error = 860.
56	93.3	Same as run no. 55 except OFFSET in SQ eliminated (for this run only).
57	94.1	Same as run no. 55 except SQ thresholds = 2,000/1,000.
58	94.1	Same as run no. 55 except maximum valley-point error = 700 and SQ thresholds = (800, 330, 860, 400, 340, 400, 820, 700, 420, 430) for digits 0 through 9, respectively.
59	93.5	Same as run no. 58 except difference data in scanning patterns eliminated.
60	93.4	Same as run no. 59 except middle column of scanning patterns eliminated.
61	94.1	Same as run no. 58 except with minor bug in subroutine DECIDE corrected.
62	94.2	Same as run no. 61 except with minor adjustments to the TE normalizing constants to account for TE changes due to T-function inclusion.
63	94.1	Same as run no. 61 except with minor changes to the quantization levels for T-function (for this run only).
64	91.9	Same as run no. 62 except T-function quantized to 16 levels maximum valley-point error = 1,100; PVR = 1.05; SQ thresholds = 2,000/1,000.

TABLE 19. SYNOPSIS OF EVALUATION RESULTS (Continued)

Run No.	Percent Correct Recognition	Remarks (Changes From Previous Runs)
65	94.9	Same as run no. 62 except TE normalizers modified to account for confusion-matrix entries from run no. 62.
66,67	95.1	
68	95.2	
69	95.3	Same as run no. 62 except TE normalizers modified to account for confusion-matrix entries from previous run.
70-72	95.2	
73	95.3	

The following additional parameters were introduced for the syntactically unconstrained digit sequence recognition algorithm added during the current contract.

Minimum (3-centisecond) and maximum (120 centisecond) interdigit times (times between first reference point of one word and last reference point of previous word)

Minimum acceptable ratio of average recognition pattern energy for each word to maximum average recognition pattern energy for all words = 0.15.

Note from table 19 that poorer results were obtained in experiments that eliminated either the floor (OFFSET) to the valley-point error (run no. 56) or the difference data from the scanning patterns (run no. 59).

The digit recognition results for each digit for selected evaluation results are shown in Table 20. The evaluation run results shown are only those exhibiting significant improvements over previous runs. These improvements occurred for run no. 52 because the total normalized error was lowered for three reference-point words, for run no. 55 because the T-function (see Section III.B for definition) quantized to 3 bits was included in the scanning patterns, and for run no. 73 because the normalizers for the recognition error were modified. (Note from Table 19 that quantizing the T-function to 4 bits in run no. 64 degraded performance.)

The change made for run no. 52 that lowered the total normalized error for longer words was a heuristic justified only by the fact that longer words (those with more reference points) are less likely to be spurious hypotheses. This same philosophy was used in the speaker-dependent recognition. These heuristic normalization constants (HNCs) used were as follows:

No. of reference points 2 3 4 5 6 7 HNC 1.00 0.90 0.81 0.73 0.66 0.59

TABLE 20. DIGIT RECOGNITION RESULTS FOR SELECTED EVALUATION RUNS

			Evaluation	Run Numbe	r	
Digit	45	50	52	55	62	73
0	91.9	91.9	94.3	94.8	95.4	95.4
1	92.6	94.0	92.9	95.2	93.1	92.4
2	76.7	78.1	88.3	92.7	92.7	94.1
3	89.4	89.3	86.7	89.3	89.2	91.9
4	88.6	88.6	88.2	89.7	90.2	97.2
5	98.5	96.2	96.1	97.6	98.1	95.9
6	95.8	95.2	97.8	98.3	98.3	98.0
7	83.7	84.4	89.9	93.8	94.4	98.0
8	97.5	96.8	93.3	97.3	98.5	93.8
9	89.8	89.1	88.0	89.3	91.2	93.7
Overall	90.5	90.3	91.8	93.8	94.2	95.3

The inclusion of the T-function in the scanning pattern was prompted by vowel/nasal reference points being moved into the nasal rather than being at the phoneme boundary. This was primarily caused by the inclusion of reference patterns to accommodate nasalized vowels. Oftentimes, even words having non-nasalized vowel-to-nasal transitions produced lower valley-point errors matching a portion of the nasalized vowel-to-nasal reference pattern. Since the desire was to favor choosing reference-point candidates at the locations of T-function peaks in the input, such a bias could be provided by including an inversely quantized value of the T-function in conveniently unused 4-bit fields in the scanning pattern for the input (Figure 36). Since the corresponding 4-bit fields of the reference scanning patterns were zero, the inverse quantization of the T-function meant that the larger T-function values (lower inversely quantized values) that usually occur at phoneme boundries would produce lower scanning errors relative to those produced when the spectral or energy change was not so great during more nearly steady-state portions of the word. The quantization thresholds given in Table 21 were derived from a cumulative distribution plot of T-function values at the selected reference points in the digit recognition design data and at ±1 ard ±2 time samples around those points.

In addition to the improved recognition performance shown for run no. 55 in Table 20, the benefit of including the T-function is the scanning pattern can also be seen by the decrease in the average recognition pattern error, indicating improved time registration of the input speech. This decrease is shown in Table 22.

The third performance improvement was prompted by the confusion matrix for evaluation run no. 62 (Table 23). Four quite large nonsymmetrical substitutions are shown in the off-diagonal entries of the confusion matrix. Since the digits are selected that minimize the minimum total normalized error across the sequence, adjustment of the relative errors among digits will affect the distribution of substitutions in the confusion matrix. The mechanism for performing this adjustment can be seen from the following equation for the total normalized error for the digit k:

$$NE_{k} = HNC_{k} \left[\frac{TE_{k}/\text{no. of column in digit } k}{TE_{k} \text{ normalizing constant}} \right] + \frac{w_{k}}{NPP} \sum_{i=1}^{NPP} PPE_{i}$$
 (29)

where HNC is the heuristic normalizing constant, TE is the recognition pattern error, NPP is the number of reference-point pairs, PPE is the point-pair error between two reference points, and w is a weighting constant for the sum of the PPEs.

The TE_k normalizing constants are calculated from the expected values of TE_k as follows:

$$TE_{k} \text{ normalizing constant} = \frac{\widehat{TE}_{k}/\text{no. columns in } k}{\frac{1}{10} \sum_{i=0}^{9} (\widehat{TE}_{i}/\text{no. columns in } i)}$$
(30)

Five sets of values for these normalizing constants are given in Table 24, derived from five different sources:

- (1) E_e + E_s + E_a from Table 20 of the Speaker Verification III report
- (2) Values of J_e for the number of reference patterns chosen in the total voice study for each digit
- (3) Values of TE for each of the digits in correct sequences in the 6-digit sequence evaluation data set for run no. 33
- (4) Same as source 3, for run no. 57, which includes the T-function in the scanning patterns
- (5) Values derived from incrementally changing the values from source 4 during run nos. 62 through 73.

Although the normalizing constants derived from source 5 will certainly give somewhat biased results since they are tuned to the evaluation set, it should be remembered that the test set is reasonably large (106 speakers). An independent test on the second set of data described in Subsection VI.A.1 showed that, while not achieving the 19-percent reduction in error rate on the six-digit sequences between run nos. 62 and 73, a 6-percent reduction in error rate was achieved using the normalizing constants from source 5 from that achieved using those from source 3.

The confusion matrix for the final evaluation run on this study (no. 73) is shown in Table 25.

One final observation made on the results of the evaluation runs was the usual problem of poorer performance for females than for males, as shown both by the overall recognition results in Table 26 and by the histogram of digit recognition performance (table 27), both from evaluation run no. 73.

3. Digit-Recognition Results for Three-Digit Sequences

Although the data used for the testing reported in the last subsection were from a large number of subjects of different ages, races, dialects, and educational backgrounds, the multiple

										DIFFERENCE	DATA							
,		t (i-2)			$\left.\right\rangle$ t (i-1)	£(i)	∑ t (i+1)	} t (i+2)	$\left.\begin{array}{c} t (i-1) - t(i-2) \end{array}\right\}$			$ \int t(i+2)-t(i+1) $		ENERGY,	ENERGY DIFF,	H-FCN		
	FILTER 4	8	12	COS REGR C2									TFCN i-2	TFCN i-1	TFCN	TFCN i +1	TFCN i+2	4 BITS
	FILTER 3	7	11	SIN REGR C1	AT AS 1-4	AT AS 1-4	AT AS 1-4	AT AS 1-4	AT AS 1-4	AT AS 1-4	AT AS 1-4	AT AS 1-4	0	0	0	0	0	——4 BITS——
	FILTER 2	ø	10	41	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	SAME FORMAT AS 1-4	ENERGY 1-2	ENERGY i-1	ENERGY	ENERGY i+1	ENERGY i+2	4 BITS
	FILTER 1	ı	б	13									E (i-1)-E (i-2)	E(i)-E(i-1)	E(i+1)-E(i)	E (i+2)-E (i+1)	0	4 BITS
MORD	-	n	m	4	5-8	9-12	13-16	17-20	21-24	25-28	29-32	33-36	37	38	39	04	14	

Figure 36. Scanning Pattern Format

TABLE 21. T-FUNCTION QUANTIZATION THRESHOLDS

Quantized Value	Range of T-Function			
0	214– ∞			
1	163-213			
2	128-162			
3	101-127			
4	79-100			
5	60-78			
6	40-59			
7	0-39			
6	40-59			

TABLE 22. DECREASE IN AVERAGE RECOGNITION PATTERN ERROR BY INCLUDING T-FUNCTION IN SCANNING PATTERNS AVERAGE RECOGNITION PATTERN ERROR

Digit	Run No. 45 (No T-Function)	Run No. 57 (3-Bit T-Function)		
0	359.8	354.0		
1	440.4	409.5		
2	316.9	311.3		
3	294.3	289.0		
4	221.5	218.8		
5	431.4	424.9		
6	283.3	275.2		
7	423.4	415.4		
8	287.3	288.3		
9	444.7	435.3		

TABLE 23. CONFUSION MATRIX FOR DIGIT RECOGNITION FOR 6-DIGIT SEQUENCES FOR RUN NO. 62

						Recogniz	zed					
		0	1	2	3	4	5	6	7	8	9	X
	0	660	8	1	10	1	2	3		1	3	
	1	4	403	1	12	6	3		1	2	2	
	2	2	1	581	21	1	-	3		17		
	3	8	1	21	657	4		~		46	1	
	4	1	13		4	644	49					
Said	5	1	3	i	1	2	779			1		8
S	6	2			4			752		5		1
	7	3	7	2	2		4		719	18	1	
	8				3	1		2		398		1
	9	2	9	1	12	1	2			11	402	1
	X		1		1			-		2		

TABLE 24. RECOGNITION ERROR (TE) NORMALIZING CONSTANTS

Digit	Source 1	Source 2	Source 3	Source 4	Source 5
-		Douite 2	Domet 5	Source 4	Source 3
0	1.039	1.022	0.998	1.005	1.020
1	1.170	1.089	1.225	1.163	1.150
2	0.708	0.837	0.878	0.884	0.890
3	1.085	1.016	1.021	1.026	1.020
4	0.853	0.719	0.681	0.691	0.738
5	1.182	1.292	1.195	1.207	1.010
6	0.705	0.870	0.788	0.782	0.730
7	1.052	0.945	0.977	0.983	1.170
8	1.007	1.133	1.008	1.024	0.860
9	1.200	1.076	1.229	1.236	1.330

TABLE 25. CONFUSION MATRIX FOR DIGIT RECOGNITION FOR 6-DIGIT SEQUENCES FOR RUN NO. 73

2

	Recognized										
3	4	5	6	7	8	9	X				
11	2	1	3	1		3	1				
11	12			1		4	1				
21	1		2	2	7	1					
678	8			1	14	1					
1	691	9		1		-					
	8	763				1	21				
7	2	-	749	1	1						

	0	657	4	1	11	2	1	3	1		3	1
	1	4	401		11	12			1		4	1
	2	2	1	589	21	1		2	2	7	1	
	3	11	1	23	678	8			1	14	1	
	4		9		1	691	9		1			
Said	5		3			8	763				1	21
S	6	2		1	7	2		749	1	1		1
	7	2	2	1	1	1	1		741	4	1	1
	8		1		4	2		7	6	380	1	4
	9	5	4	1	10		1			5	413	2
	X		1	1	1						2	

TABLE 26. DIGIT-RECOGNITION PERFORMANCE OF MALES AND FEMALES

Percent Correct Recognition

Digit	Males	Female
0	97.8	91.6
1	92.8	91.8
2	94.4	93.6
2 3	94.0	88.2
4 5	99.8	93.1
5	97.0	94.1
6	99.1	96.7
7	99.2	96.1
8	94.9	92.4
9	95.5	91.4
Average	96.7	93.1

TABLE 27. HISTOGRAM OF DIGIT-RECOGNITION PERFORMANCE

No. of Errors	Percent Correct	Number of Subjects	Number of Males	Number of Females
0	100	22	19	3
1	98	20	14	6
2	97	21	13	8
3	95	11	6	5
4	93	11	7	4
5	92	4	1	3
6	90	5	0	5
7	88	4	2	2
8	87	3	1	2
9	85	0	0	0
10	83	3	0	3
11	82	0	0	- 0
12	80	0	0	0
13	78	i	1	0
14	77	0	0	0
15	75	i	0	1
≥16	<75	0	0	0

evaluations using these same data made a further independent test mandatory. Results in this section use the three-digit sequences of the second data base, excluding those with the "oh" pronunciations. Results in the next subsection are for the digits said in isolation from the same data base.

Note that, for these 50 sequences, all digits appear an equal number of times and in all contexts of preceding and following digits. Since the original application of the work done in the total voice study was for syntactically constrained sequences, not all digit pairs were used in the design data, as described in the previous subsection. Hence, the recognition performance is expected to be poorer for digits involved in these transitions. This poorer performance does not, however, reflect on the method developed for choosing reference patterns, but only on the inadequacy of the design data for unconstrained digit recognition.

The data used in these tests were collected in sound booths or sound-treated rooms at two locations: Texas Instruments (Dallas, Texas) and the Institute for Advanced Study of the Communication Process (Gainesville, Florida).

Since the poorer performance of females has been demonstrated in the previous subsection (as well as in all other word-recognition studies that have been done, to the best of the author's knowledge), experiments were performed for male subjects only, 12 from Dallas and 11 from Gainesville. The recognition performance is shown as the far right column in the confusion matrix in Table 28. The overall percent correct is 94.0.

It has been noted in these studies that the confusion matrix can be quite speaker-dependent. Since, in the 3-digit sequence data base, there were more digits (150 versus 60) and fewer subjects (23 versus 106), the confusion matrix entries are more susceptible to high substitution rates by particular speakers. For example, 23 of the 48 3-for-2 substitutions shown in Table 28 were caused by two speakers. However, the two digits with the greatest reduction in recognition rate (2 and 8) from that given in Table 26 for males are two of the four digits

TABLE 28. CONFUSION MATRIX FOR DIGIT RECOGNITION FOR 3-DIGIT SEQUENCES CONSTRAINED IN LENGTH

						Recogn	nized					Percent
	0	1	2	3	4	5	6	7	8	9	X	Correct
0	340	1								3		98.8
1	1	331		2	3					8		95.9
2	7	1	287	48						1	1	83.2
3	3	3	5	327					1	6		94.8
4	1	5		1	330	6					1	95.9
Said 2						341				1	3	98.8
6	14		1	1	1		323	3	1			94.0
7		1	1	1				341	1			98.8
8		2	1	29	1		5	2	294	10	1	85.2
9		11			1				6	326	1	94.5
X				1			1			2		

having reference points located at the word boundaries. Hence, if the reference scanning patterns used for locating these points do not explicitly account for all allowable contexts, then these reference points may be missed during scanning because of the lack of a significantly deep valley in the scanning error (distance to the reference scanning pattern). In such a case, even though the time-normalized recognition pattern is much less affected by context, this digit would not even be hypothesized because of the missing reference point.

Since the percent correct achieved for the 3-digit connected digits was 94.0 using reference patterns that did not account for all transitions, it is reasonable to assume that the 96.7-percent correct achieved for males in the 6-digit sequences was minimally, if any at all, caused by tuning to the test set during the parameter evaluations done on the 6-digit sequence test set. It appears reasonable that such a recognition rate could be achieved on the unconstrained digit recognition if the reference scanning pattern set were expanded to include patterns to account for these transitions. Such patterns could be generated with an expanded design set using the clustering techniques developed in these studies.

4. Digit Recognition Results for Isolated Digits

A very limited experiment was run to test the digit recognition on isolated digits. The test involved two samples of each of the 10 digits said in isolation from the same 23 speakers used in the 3-digit sequence test. In view of the recognition rates achieved on continuous speech, the recognition rate of the isolated digits was a surprisingly low 95.4 percent, with over one-third of the errors being caused by 3-for-2 and 5-for-4 substitutions. More expected would be results such as that achieved by Martin², whose error rate was approximately halved between his test T-2 (Philadelphia—connected digits) and his test T-3 (Philadelphia—isolated digits).

One possible explanation for the higher-than-expected isolated digit error rate is that the patterns contained in the reference set to account for contextual variations are generating spurious hypotheses during isolated word recognition.

The specific confusion matrix for the isolated digits is given in Table 29.

TABLE 29. CONFUSION MATRIX FOR DIGIT RECOGNITION OF ISOLATED DIGITS

						Reco	gnized					
		0	1	2	3	4	5	6	7	8	9	Percent Correct
	0	45			1							97.8
	1		45								1	97.8
	2			42	4							91.3
	3				46							100.0
Said	4		2			40	4					87.0
Š	5						46					100.0
	6				1			45				97.8
	7	-						-	44		2	95.7
	8				3					43		93.5
	9						1			2	43	93.5

5. Effect of Spectral Normalization Technique on Digit-Recognition Performance

Subsequent to the performance tests described previously in this section, investigation of the high 3-for-2 substitution rate in the 3-digit sequences revealed a mechanism for improving recognition results and for making them less susceptible to variations in background noise. This investigation revealed that the valley point error for reference point 1 for the digit 2 (i.e., the silence/plosive transition) had a higher error for the real-time data than for the digitized data. This was found to be caused, at least in part, by different "silence" spectra. The differences were alleviated to some extent by increasing σ_{\min} in the constant σ_j^* used in normalizing the regressed filter outputs (see discussion in Appendix A):

$$\sigma_{j}^{*} = \sigma_{\text{post}_{i}} + \sigma_{\min} \tag{31}$$

where $\sigma_{\mathrm{post}_{\hat{\mathbf{i}}}}$ is the post-regression standard deviation.

The specific effects on the spectrum of changing σ_{min} can be seen from the three spectra shown in Figure 37 for the word two as said by J.S. The increase in the number of hypothesized digits and the decrease in the total normalized errors for the six digits in the sequences from which the spectra in Figure 37 were extracted are shown in Table 30. A further breakdown into the terms that make up the total normalized errors is given in Table 31. Since these tests were done directly from the analog tape for each trial, exactly repeatable filter outputs are not obtained. However, general trends can be noted from these tables, such as the over 50-percent drop in the valley point error for reference point 1 of the digits 5, 2, and 4, which contain silence or low energy frication (/f/). An additional point of interest is the consistently lower recognition errors for the digitized data than for the analog data.

The expected benefits to be derived from using a larger σ_{\min} are twofold. First, the normalized spectrum will tend to be more even for silence or low-energy fricatives, making the resulting patterns more resistant to variations in background noise. Second, since the reference

```
TO THE TOTAL TOTAL
                                                                       \sigma_{\min} = 62
 (B)
                                                                                                                                                                                                                                               TITLE TITLE THE TITLE TO THE TITLE THE TITLE TO THE TITLE
                                                                                                                                                                                                                        \sigma_{\min} = 250
(C)
                                                                                                                                                                                                                   TITITE TO THE TOTAL TOTAL
```

Figure 37. Spectra for Digit "Two" for Speaker J.S.

TABLE 30. TOTAL NORMALIZED ERROR (NE) FOR DIGITS FROM SEQUENCE 852–734 FOR SPEAKER J.S.

	No of Hypothesized			Normaliz	ed Error		
σ_{\min}	Digits	8	5	2*	7*	3	4
18	49	55	55	63(70)	48(54)	59	66
37	58	60	49	49(55)	47(53)	64	57
62	52	54	44	45(51)	50(56)	52	55
100	58	51	41	45(50)	46(52)	58	51
150	76	51	38	40(45)	48(54)	52	46
200	70	40	39	42(47)	45(51)	50	45
250	84	38	38	41(46)	45(51)	47	45
375	72	40	33	36(41)	45(50)	43	44
Run No. 53	~52	47	38	40(45)	44(49)	46	54
$(\sigma_{\min} = 62)$							

^{*}Since three reference point digits are multiplied by 0.9, the unadjusted errors are given in parentheses.

TABLE 31. VALLEY POINT ERRORS, SEQUENCE ERRORS (SQ), AND RECOGNITION ERRORS (TE) FOR DIGITS FROM SEQUENCE 852–734 FOR SPEAKER J.S.

			8				5				2		
σ_{\min}	1	2	SQ	TE	1	2	SQ	TE	1	2	3	SQ	TE
18	389	161	256	335	385	225	312	468	442	330	320	657	418
37	346	149	216	395	265	218	252	430	308	253	293	474	340
62	338	137	202	358	264	200	214	389	291	261	268	447	320
100	250	124	152	346	279	167	198	373	255	287	250	404	328
150	324	130	190	331	181	188	158	359	206	244	231	326	297
200	240	105	136	268	214	165	164	368	213	271	231	367	308
250	210	111	126	260	218	163	166	362	181	257	273	358	298
375	212	103	122	275	191	120	124	318	195	247	220	313	273
Run No. 53	303	143	190	309	222	219	200	343	144	261	253	305	308
			7					3				4	
σ_{min}	1	2	7	SQ	TE	1	2	3 SQ	TE	1	2	4 SQ	TE
σ _{min}	1 280	2 272		SQ 353	TE 489	1 406	2 352		TE 306	1 376	2 129		TE 336
			3					SQ				SQ	
18	280	272	3 148	353	489	406	352	SQ 448	306	376	129	SQ 214	336
18 37	280 336	272 267	3 148 153	353 392	489 468	406 420	352 278	SQ 448 404	306 358	376 239	129 129	SQ 214 152	336 302
18 37 62	280 336 352	272 267 263	3 148 153 220	353 392 487	489 468 471	406 420 374	352 278 320	SQ 448 404 388	306 358 266	376 239 193	129 129 114	SQ 214 152 122	336 302 298
18 37 62 100	280 336 352 318	272 267 263 255	3 148 153 220 143	353 392 487 382	489 468 471 462	406 420 374 350	352 278 320 271	SQ 448 404 388 356	306 358 266 329	376 239 193 202	129 129 114 108	SQ 214 152 122 122	336 302 298 278
18 37 62 100 150	280 336 352 318 300	272 267 263 255 259	3 148 153 220 143 125	353 392 487 382 359	489 468 471 462 488	406 420 374 350 314	352 278 320 271 251	SQ 448 404 388 356 310	306 358 266 329 296	376 239 193 202 181	129 129 114 108 88	SQ 214 152 122 122 102	336 302 298 278 251
18 37 62 100 150 200	280 336 352 318 300 291	272 267 263 255 259 247	3 148 153 220 143 125 122	353 392 487 382 359 324	489 468 471 462 488 464	406 420 374 350 314 306	352 278 320 271 251 238	SQ 448 404 388 356 310 294	306 358 266 329 296 290	376 239 193 202 181 189	129 129 114 108 88 92	SQ 214 152 122 122 102 108	336 302 298 278 251 248

patterns used in speaker-independent digit recognition result from averaging patterns from many speakers, the reference patterns appear more "washed-out," lacking the sharp contrasts found in speaker-specific patterns. Hence, a flatter spectrum on the input speech resulting from using a larger σ_{\min} would probably result in a better match to speaker-independent reference patterns.

This hypothesis was tested using a tape generated at RADC (in the computer room containing speech-processing equipment) consisting of two repetitions each of two speakers (R.V. and J.F.) of the 50 three-digit sequences given in Table 18. The recognition results for all four repetitions are given in Table 32 using both a σ_{\min} of 62 and a σ_{\min} of 250. The recognition results for the larger σ_{\min} show a small improvement. A performance improvement would correspondingly be expected on the results presented in the previous subsections since all these results used data preprocessed using a σ_{\min} of 62.

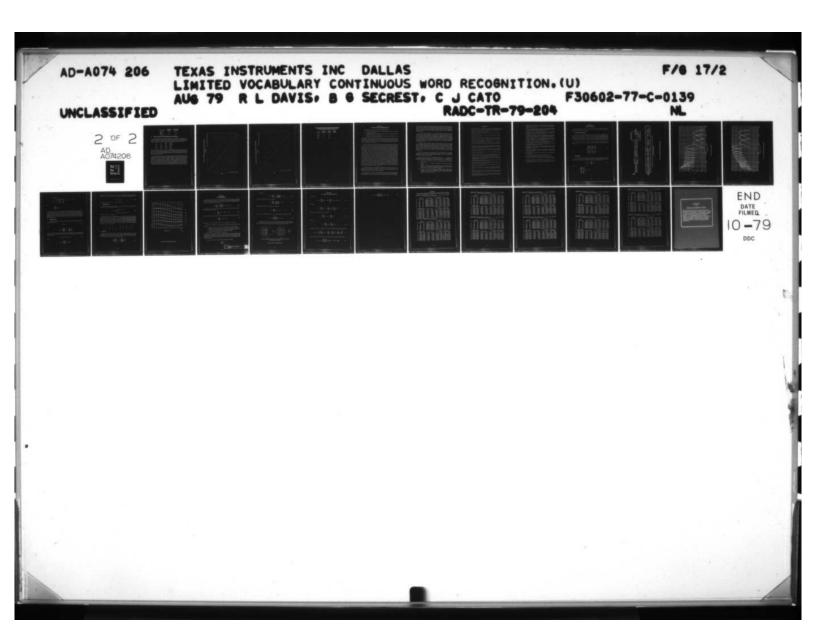
TABLE 32. σ_{min} VARIATION PERFORMANCE TEST

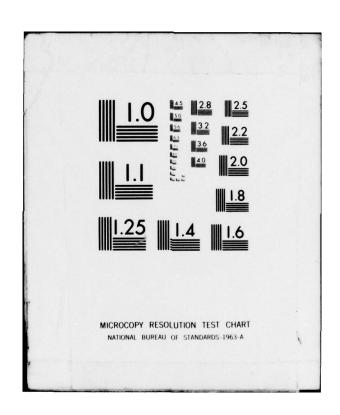
		Length Co	onstrained	Length Unconstrained					
Subject	Session	$\sigma_{\min} = 62$	$\sigma_{\min} = 250$	$\sigma_{\min} = 62$	$\sigma_{\min} = 250$				
J.F.	1	96.0	93.3	98.0	95.3				
J.F.	2	95.3	97.3	96.7	96.7				
R.V.	1	89.3	90.0	90.0	92.0				
R.V.	2	91.3	94.0	92.0	92.0				
Overall		92.5	93.7	94.2	94.0				

In conclusion, an observation made concerning sequence recognition during the testing of these two subjects should be noted. In an operational system, where sequences can be repeated, the only consequence of rejected sequences is a decrease in throughput (assuming a sequence can finally be accepted if repeated); therefore, rejections should not be used in calculation of the percent correct recognition rate. The percent correct sequence recognition rate is then given by

$$\% \text{ correct} = \frac{\text{no. of correct sequences}}{\text{no. of correct sequences}}$$
or
$$\% \text{ correct} = \frac{\text{no. of correct sequences}}{\text{no. of correct sequences}}$$

In the case where the length is constrained, the present tree-searching algorithm described in Section II chooses the best sequence of the specified length. However, the sequence recognition results for the limited testing given in this subsection indicate that by modifying the algorithm so that the sequence is accepted only if the length of the best sequence is the same as the specified length, the sequence recognition rate (as defined above) would improve as shown in the following table:





No. of Utterances	Best Sequence of Length 3	If Best Sequence Is of Length 3
Correct	336	330
Rejected	7	25
Incorrect	57	45
Percent correct	85.5	88.0

B. LIMITED VOCABULARY WORD-RECOGNITION EXPERIMENT

Five speakers from the Speech Research Branch at Texas Instruments were tested on the limited word-recognition algorithm using an automatic enrollment and a hand enrollment. Each speaker was recorded onto analog tape while seated in the sound booth. An enrollment session collected from each speaker consisted of four discrete repetitions of the following words:

Zero	Five	Minus	Hundred
One	Six	Plus	Thousand
Two	Seven	Point	Enter
Three	Eight	Backup	Erase
Four	Nine	Punch	Display

At some time later the same day or the next day, an execution session was collected from each speaker. The execution session consisted of a set of 20 phrases of three randomly chosen words and a set of 20 phrases of random lengths (up to seven words) of randomly chosen words. Each phrase in the execution session was spoken continuously. Two of the speakers had two execution sessions spaced a half day apart.

The speakers were then enrolled off-line using both automatic enrollment and hand enrollment. The execution sessions were then tested against these enrollments. The results of the experiment are given by the confusion matrices of Tables 33 and 34. The left of the matrix shows what was said and the top of the matrix shows what was recognized. An entry in the "X" column means nothing was recognized (a deletion). The entries in the matrix are the number of times a word was recognized versus what was said. A compilation of the results for each inidividual speaker is given in Table 35. One of the speakers (Keith) had two execution sessions, and his first execution was used for a supervised updating. The second execution session was used against the updated reference patterns. The results are given in the last column of Table 35.

TABLE 33. CONFUSION MATRIX FOR AUTOMATIC ENROLLMENT

	×	4	7	28	6	=	12	4	2	2	00	e	7	7	7	=	7	-	=	-	-
	DIS	1	1	1	- 1	-	1	1	1	1	1	1	-	1	1	1	1	1	1	1	4
	ER	1	1	1	-	1	1	1	1	1	1	1	1	-	1	1	1	1	1	84	1
	EN	1	-	-	1	-	1	1	1	1	1	1	1	1	1	1	1	1	33	1	1
	E	-	1	-	-	-	-	1	1	1	1	1	1	-	1	1	1	31	1	1	1
	H	-	1	-	1	1	1	1	1	1	1	1	1	1	1	-	42	1	7	-	1
	P	1	1	1	-	1	1	1	1	1	1	1	-	1	1	25	1	1	1	1	1
	BU	1	1	1	1	1	1	1	1	1	1	1	1	-	51	1	1	-	1	1	1
		1	7	1	7	-	1	1	1	1	1	1	1	30	1	1	-	1	1	1	-
Recognized	+	1	9	1	-	-	8	1	7	-	1	1	45	1	1	-	1	-	1	1	1
Recoi	1	1	1	1	-	1	-	1	1	1	-	38	1	1	1	1	1	1	1	1	1
	6	1	1	1	1	1	1	1	-	1	45	1	1	1	1	1	1	1	1	1	1
	∞	1	1	1	7	1	1	1	1	34	1	1	1	1	1	1	1	1	1	7	1
	-	1	1	1	1	1	1	1	4	1	1	-	1	1	1	1	1	1	1	1	1
	9	7	1	1	-	1	1	43	7	-	1	1	1	-	1	1	7	1	-	1	1
	S	-	-	-	-	1	23	-	1	-		1	1	1	1	-	1	1	1	1	1
	4	1	-	-	1	32	-	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	3	1	-	7	32	1	1	1	1	1	1	1	1	-	1	1	-	1	1	-	1
	7	1	1	35	7	1	1	1	1	-	1	1	1	1	1	1	1	-	1	1	-
	-	1	53	1	1	1	1	1	1	1	1	1	1	1	1	7	1	1	1	1	1
	0	36	1	1	-	1	1	1	1	1	1	1	1	1	1	1	1	-	1	1	1
		0	_	7	8	4	S	9	1	00	6	1	+		BU	P	HI	H	EN	ER	DIS
										P	Sai										

TABLE 34. CONFUSION MATRIX FOR AUTOMATIC ENROLLMENT

	×	c	7	16	10	12	4	3	-	2	1	S	4	-	-	w	4	7	12	1	1
	DIS	1	1	1	1	-	1	1	1	1	1	1	1	1	1	1	-	1	1	1	46
	ER	1	1	1	1	1	1	}	1	1	1	1	-	1	-	1	1	-	-	53	1
	EN	-	1	4	-	1	1	-	1	1	1	1	1	1	1	1	1	1	53	1	1
	Ŧ	1	1	1	1	1	1	-	1	1	-	-	-	1	1	1	1	34	1	1	1
	H	1	1	1	1	1	1	1	1	1	1	-	1	-	1]	20	1	1	1	-
	P	1	1	1	1	1	1	1	-	1	-	1	1	1	-	34	1	1	1	1	1
	BU	1	1	-	-	-	1	1	1	1	-	1	1	1	52	1	1	- 1	-	1	1
		.1	4	1	1	-	1	1	1	-	1	1	1	37	1	n	1	1	1	1	1
_	+	1	1	1	1	-	1	1	1	-	1	-1	4	-	1	1	1	1	1	1	1
Recognized	1	1		1	-	1	1	1	1		1	39	- 1	1	1	1	-	-	-	-	1
Reco	6	-1	1	1	1	1	1	-	1	1	51	1	1	1	1	1	1	1	1	1	1
	œ	1	-	m	7	1	-	-	1	37	1	1	1	1	1	1	1	1	1	-	1
	~	1	1	1	1	1	1	1	36	1	1	1	1	1	}	1	1	1	1	-	1
	9	1	1		-	-	-	3	1	-	1	-	1	-	-	1	1	-	1	-	1
	S	i	-	1	-	1	39	-	-		1	1	1	1	-	1	1	1	1	1	1
	4	1	1	-	1	33	1		-	1	1	1	1	-	1	- 1	1	1	1	1	1
	8	-	1	-	35	1	1	1	1	1	-	1	1	1	1	-	-	1	-	-	1
	7	-	1	4	-	-	-		1	1	-	- 1	1	- 1	1	1	1	1	1	1	1
	-		36	-	1	1	-	-	1	-	-	7	1	-	1	-	1	1	1	1	1
	0	35	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	-	1	1	1
		0	-	7	6	4	S	9	1	∞	6	1	+		BU	PU	HU	IH	EN	ER	DIS

TABLE 35. PERCENT CORRECT SPEAKER-DEPENDENT RECOGNITION RESULTS FOR CONTINUOUS UTTERANCES FROM A 21-WORD VOCABULARY ENROLLED ON ISOLATED WORDS

Speaker	Automatic Enrollment	Hand Enrollment	Updating (Hand)
Gene	77	78	
Richard	73	88	
Louise	74	87	
George	76	94	
Keith	66	80	90

SECTION VII

CONCLUSIONS AND RECOMMENDATIONS

The three major areas of research during this study contract were

- (1) High-performance, speaker-independent, connected-digit recognition for syntactically unconstrained digit sequences
- (2) Clustering algorithms for use in the development of sets of reference patterns for speaker-independent word recognition
- (3) Automatic enrollment for speaker-dependent, connected-word recognition for syntactically unconstrained word sequences.

The program culminated in the installation of the speaker-independent, connected-digit recognition program on the BISS-ADM speaker verification system at RADC using the total voice reference patterns for compatibility. In addition, a long-standing hardware failure with the digital filters on the BISS-ADM system was corrected, resolving performance discrepancies between the systems at RADC and Texas Instruments.

As part of the three tasks, several developments resulted that are generally applicable to the speech technology used in this study. The first of these is a modification to the algorithm for searching the table of hypothesized words (directed graph) that significantly reduces the processing time. The second development is a technique (transparent to previous programs) for including a measure of the spectral transitionitivity (T-function) in the scanning patterns for the purpose of improving the time registration of reference-point locations. The third development is the capability of digitizing and playing back speech data through A/D and D/A connections to the fast array processor. This provides the basis for the fourth development, which is simulation of the digital filters in the array processor, allowing parametric variation of the filter-bank definition and the consequent ability to perform a variety of tests with data that can be more precisely replicated using a variety of filter-bank definitions. The new speech channel capability was also necessary for a fifth development, that of using a quantized autocorrelation value out of an autocorrelation pitch tracker previously implemented on the array processor to produce a "soft" voicing decision for each frame of filtered speech data for eventual incorporation into the time-normalized recognition pattern. The sixth general development came as a natural extension to the capabilities provided by the speech channel and filter simulation. This development is the programming of the preprocessing function in the array processor and subsequent amalgamation of digitizing, filtering, and preprocessing in the array processor for inputting preprocessed speech data to the word-recognition programs. This capability reduces the 980B processing time by about 35 percent, allowing the word recognition algorithm that uses the new directed-graph searching algorithm to operate sufficiently fast to allow continuous speech input without having to discontinue sampling after the input of an utterance.

The speaker-independent, connected-digit recognition portion of this study resulted in a significantly faster algorithm with a 50-percent decrease in error rate over the course of this study—from 90.5 percent correct recognition to 95.3 percent on an evaluation data set of ten 6-digit sequences from 106 speakers (64 males, 42 females).

The development of the clustering algorithm resulted in a two-stage, four-path algorithm with the mechanisms for detecting outlying data points in the design data and with subsequent analysis routines for comparing the results from the various paths and testing the validity of

resulting clusters on the basis of comparisons with a priori information about the design data set. The results of the analysis of the digit-recognition design data set revealed that, although the clusters selected during the total voice speaker verification contract generally were good partitions of the data, use of partitions resulting from other paths in the more comprehensive algorithm would have resulted in somewhat more compact clusters in terms of minimizing the sum-of-squared error.

The research into development of an automatic enrollment technique for speaker-dependent word recognition resulted in a method that yielded very good results for isolated word recognition but less acceptable results when used in continuous speech from the same speaker. The better results achieved with comparable hand enrollments point to the desirability of a semiautomated enrollment procedure allowing the operator the option of modifying reference-point locations and recognition-pattern format definitions defined by an automated front end. Independent of the enrollment method, however, the benefit of reference file updating as a means of accommodating contextual variability, as well as intersession variability, became abundantly clear.

Throughout all three phases of this study, the general limitation existed of an insufficient speech data sample rate and spectral resolution of the filter bank, especially in the higher frequency bands. This limitation must be removed before any further word recognition development. In addition, although all recognition features up to this point have been spectral amplitudes or direct correlates thereof (regression coefficients and energy), it is time that more features are used. This, in fact, was the impetus behind the addition of the "soft" voicing decision (quantized autocorrelation coefficient) to the spectral parameters derived during preprocessing.

Care must be taken, however, that none of the new features added are subject to measurement errors sufficient to actually degrade performance. In addition to not degrading overall performance, new features must not degrade performance of the poor speakers while improving the results for the good speakers.

However, new features such as autocorrelation values or formant values will require computation capabilities exceeding those of a 16-bit minicomputer. The recommendation for future word-recognition development is that such research be done with a computing facility that

- (1) Is capable of fast arithmetic both for longer word-length integers and for floatingpoint numbers
- (2) Contains a large (1/4 to 1/2 million words) primary storage with virtual memory capability
- (3) Contains an operating system with more programmer-directed features than typically available on 16-bit minicomputers, allowing more of the time now spent on program development to be spent on speech algorithm development
- (4) Contains a fast array processor capable of performing the filter simulation, linear predictive coefficient (LPC) computations, formant tracking, autocorrelation computation, etc., necessary for the extended feature set that is required for further recognition performance improvement.

REFERENCES

- R.L. Davis, B.M. Hydrick, and G.R. Doddington, "Total Voice Speaker Verification," Rome Air Development Center Technical Report, RADC-TR-78-260, January 1978, A065160.
- T.B. Martin, "Acoustic Recognition of a Limited Vocabulary in Continuous Speech," Ph.D. Dissertation, University of Pennsylvania, 1970.
- D.R. Reddy, "Speech Recognition by Machine: A Review," Proceedings of the IEEE, 64:501-531, April 1976.
- G.R. Doddington, "A Method of Speaker Verification," Ph.D. Dissertation (Thesis), University of Wisconsin, 1970.
- H. Sakoe and S. Chiba, "A Dynamic Programming Approach to Continuous Speech Recognition," Proceedings of the 7th International Congress on Acoustics, August 1971.
- V.M. Velichko and N.G. Zagoruiko, "Automatic Recognition of 200 Words," International Journal Man-Machine Studies, 2:223, June 1970.
- F. Itakura, "Minimum Prediction Residual Principle Applied to Speech Recognition," IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-23:67-72, February 1975.
- 8. G.M. White and R.B. Neely, "Speech Recognition Experiments With Linear Prediction, Bandpass Filtering, and Dynamic Programming," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-24:183–188, April 1976.
- 9. H. Sakoe and S. Chiba, "Dynamic Programming Algorithm Optimization for Spoken Word Recognition," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-26:43-49, February 1978.
- B.T. Lowrre, "The HARPY Speech Recognition System," Ph.D. Dissertation (Thesis), Carnegie-Mellon University, 1976.
- G.M. White, "Continuous Speech Recognition: Dynamic Programming, Knowledge Nets and HARPY," Paper 28-2, 1978 WESCON Professional Program, September 1978.
- 12. J.E. Porter, "LISTEN: A System for Recognizing Connected Speech Over Small, Fixed Vocabularies, in Real Time," Naval Training Equipment Center Technical Report, NAVTRAEQUIPCEN 77-C-0096-1, April 1978.
- G.R. Doddington, "Speaker Verification," Rome Air Development Center Technical Report, RADC-TR-74-179, April 1974. 785135/5GI.
- J.S. Kenyon and T.A. Knott, A Pronouncing Dictionary of American English, G. & C. Merriam Company (Springfield, Massachusetts, 1953).
- F. Jelinek, "Continuous Speech Recognition by Statistical Methods," Proceedings of the IEEE, 64:532-556, April 1976.
- N.R. Dixon and H.F. Silverman, "The 1976 Modular Acoustic Processor (MAP)," IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-25:367-379, October 1977.
- F. Jelinek, L.R. Bahl, and R.L. Mercer, "Design of a Linguistic Statistical Decoder for the Recognition of Continuous Speech," *IEEE Transactions on Information Theory*, IT-21:250-256, May 1975.
- L.R. Bahl and F. Jelinek, "Decoding for Channels With Insertions, Deletions, and Substitutions With Applications to Speech Recognition," *IEEE Transactions on Information Theory*, IT-21:404-411, July 1975.

- L.R. Rabiner, "On Creating Reference Templates for Speaker Independent Recognition of Isolated Words," IEEE Transactions on Acoustics, Speech and Signal Processing, ASSP-26:34-42, February 1978.
- S.E. Levinson et al., "Interactive Clustering Techniques for Selecting Speaker-Independent Reference Templates for Isolated Word Recognition," *IEEE Transactions on Acoustics,* Speech and Signal Processing, ASSP-27:134-141, April 1979.
- L.R. Rabiner et al., "Speaker-Independent Recognition of Isolated Words Using Clustering Techniques," Proceedings of the International Conference on Acoustics, Speech and Signal Processing, 574-577, Washington, D.C., 2-4 April 1979.
- 22. L.R. Rabiner and J.G. Wilpon, "Considerations in Applying Clustering Techniques to Speaker-Independent Word Recognition," *Proceedings of the International Conference on Acoustics, Speech and Signal Processing*, 578–581, Washington, D.C., 2–4 April 1979.
- 23. K. Tanaka, "A Standard Category Pattern-Making Method With Application to Phoneme Recognition," *Proceedings of the Fourth International Joint Conference on Pattern Recognition*, 1030–1032, Kyoto, Japan, 7–10 November 1978.
- K. Tanaka, "A Talker Clustering Method for Standard Pattern Making," Progress Report on Speech Research '77, Electrotechnical Laboratory, Japan, August 1978.
- B. Gold, "Word-Recognition Computer Program," Massachusetts Institute of Technology, Cambridge, RLE Technical Report 452, June 1966.
- L.R. Rabiner et al., "Considerations in Dynamic Time Warping Algorithms for Discrete Word Recognition," *IEEE Transactions on Acoustics, Speech and Signal Processing*, ASSP-26:575-582, December 1978.
- K. Fukunaga and L.D. Hostetler, "The Estimation of the Gradient of a Density Function, With Applications in Pattern Recognition," *IEEE Transactions on Information Theory*, IT-21:32-40, January 1975.
- 28. R.O. Duda and P.E. Hart, *Pattern Classification and Scene Analysis*, John Wiley and Sons (New York, 1973).
- 29. J.A. Hartigan, Clustering Algorithms, John Wiley and Sons (New York, 1975).
- 30. B. Everitt, Cluster Analysis, Heinemann Educational Books, Ltd. (London, 1974).
- 31. T. Calinski and J. Harabasz, "A Dendrite Method for Cluster Analysis" (unpublished), 1971.
- 32. M.R. Anderberg, Cluster Analysis for Applications, Academic Press (New York, 1973).
- 33. F.M. Reza, An Introduction to Information Theory, McGraw-Hill Book Company (New York, 1961).

APPENDIX A SPEECH PROCESSING

The speech processing used in this study is based on the relative spectrum of speech as a function of time, which is the output of a 16-channel digital filter bank that has been preprocessed as described in this appendix.

1. FILTER BANK DEFINITION

The spectrum is obtained by processing the speech signal through a digital filter bank preceded by a first-order differencing network (for preemphasis). The filter bank consists of 16 bandpass filters, each followed by a full-wave rectifier and a four-pole lowpass Bessel filter with a 3-dB cutoff at 30 Hz. Each of the 16 filters is sampled 100 times per second. A block diagram of the spectral analysis hardware is shown in Figure A-1. Actual filter responses appear in Figure A-2 for the bandpass filters alone and in Figure A-3 for the bandpass filters with preemphasis.

For processing, the top three filters are summed and filter 14 is replaced by this sum. Filters 15 and 16 are set to zero. The resulting 14 filter outputs at each time sample are represented by:

$$A_{j} = \begin{cases} a_{1j} \\ a_{2j} \\ \vdots \\ a_{14j} \end{cases} = \begin{cases} a_{1}(t_{j}) \\ a_{2}(t_{j}) \\ \vdots \\ a_{14}(t_{j}) \end{cases}$$

$$(A-1)$$

2. REGRESSION

It has been found that, by eliminating the gross aspects of the spectrum, such as the slope and curvature, more clearly defined formant frequencies are obtained. Therefore, the spectral amplitude vector is regressed by the first three elements of an orthonormal basis set:

$$(\vec{A}_j)_R = \vec{A}_j - \sum_{k=0}^2 c_{jk} \vec{F}_k$$
 (A-2)

where

$$\vec{F}_{k} = \begin{cases} f_{1k} \\ \vdots \\ f_{14k} \end{cases} k = \{0, 1, 2\}$$

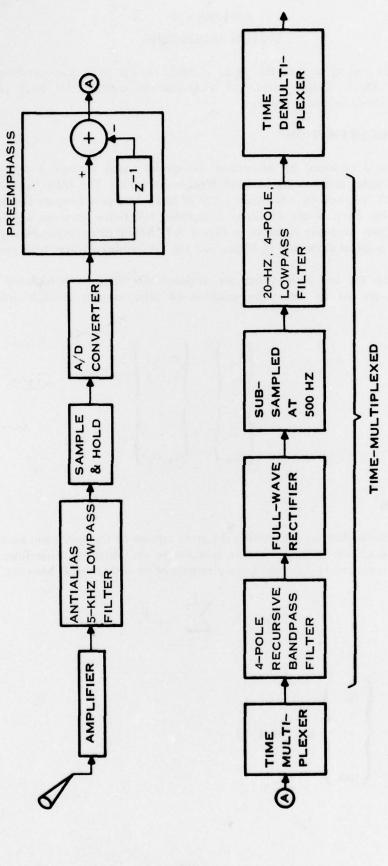


Figure A-1. Spectral Preprocessing Functional Block Diagram

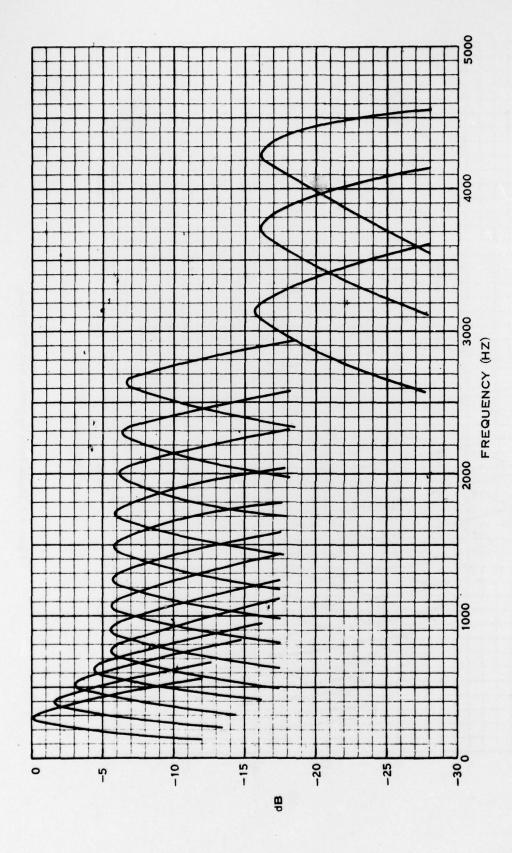


Figure A-2. Digital Filter Responses

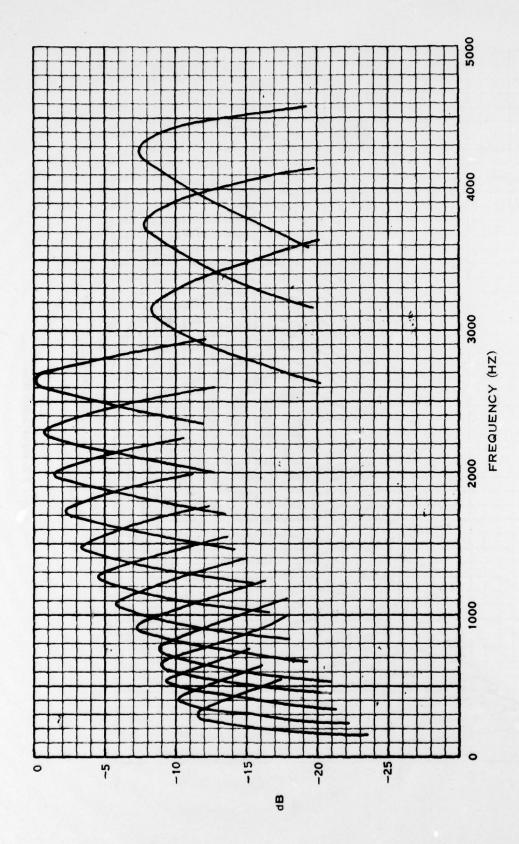


Figure A-3. Digital Filter Responses With Preemphasis

$$f_{i0} = \frac{1}{\sqrt{14}}$$

$$f_{i1} = -\frac{1}{\sqrt{7}} \sin \left[\frac{(i - \frac{1}{2})}{14} \pi \right] \qquad \{i = 1, 2, \dots, 14\}$$

$$f_{i2} = -\frac{1}{\sqrt{7}} \cos \left[\frac{(i - \frac{1}{2})}{14} \pi \right]$$

and

$$c_{jk} = \sum_{m=1}^{14} a_{mj} f_{mk}$$

Thus, the regression tends to flatten the spectrum, removing any half-cycle sine or cosine wave trends of the spectrum at time t_j . An example of a spectral waveform having a large positive c_1 is a nasal, which has one peak near the low end and one near the high end of the spectrum (around 250 Hz and 2200 Hz). An example of a spectral waveform with a large positive c_2 is a sibilant, having most of its energy above 3000 Hz. Most vowels, however, have the opposite spectral tilt because of the glottal source spectral decay with increasing frequency, yielding a large negative value of c_2 .

3. NORMALIZATION

The regressed amplitude vector is next normalized by a modified postregression standard deviation, σ_i^* for time t_i :

$$\sigma_j^* = \sigma_{\text{post}_j} + \sigma_{\text{min}} \tag{A-3}$$

where

$$\sigma_{\text{post }j}^2 = \frac{1}{11} \left(\sum_{m=1}^{14} a_{mj}^2 - \sum_{k=0}^2 c_{jk}^2 \right)$$

and $\sigma_{min} = 62$ for this study. However, it has been noticed that regression sometimes eliminates too much of the variance of the filter output vector A_j . To limit the regression, a limit is placed on σ_{post} as follows:

$$\sigma_{\text{post}_j} = \max \left(\sigma_{\text{post}_j}^2, R_{\text{min}}^2 \sigma_{\text{pre}_j}^2 \right)$$
 (A-4)

where

$$\sigma_{\text{pre }j}^2 = \frac{1}{13} \left(\sum_{m=1}^{14} a_{mj}^2 - c_{j0}^2 \right)$$

and $R_{min} = 0.6$. Note that, when $\sigma_{post_j} = R_{min} \sigma_{pre_j}$, the regression coefficients c_1 and c_2 are reduced in order to decrease the amount of regression. The resulting normalized amplitude vector is:

$$(\vec{A}_j)_N = \frac{1}{\sigma_j^*} (\vec{A}_j)_R \tag{A-5}$$

The regression coefficients c_1 and c_2 are also normalized by σ_i^* .

4. QUANTIZATION

The regressed and normalized amplitude vector is then quantized to one of eight levels according to a set of quantization thresholds ϕ_{iq} :

$$(a_{ij})_Q = q \text{ IFF}$$

$$\begin{cases} (a_{ij})_N \ge \phi_{iq} \\ (a_{ij})_N < \phi_{i,q+1} \text{ for } q = 0, 1, ..., 7 \end{cases}$$
 (A-6)

where $\phi_{iq} < \phi_{i,q+1}$; $\phi_{i0} = -\infty$; and $\phi_{i8} = \infty$.

Rather than have these quantization levels (ϕ_{iq}) being chosen to yield a uniform probability, however, it was more desirable to have the quantization thresholds cluster at higher energy levels. In this way, the sensitivity to noise can be reduced and quantization resolution is increased in the region of interest (which is the spectrum amplitude at the formant frequencies). The actual procedure used to determine the quantization thresholds is described in more detail in the total voice verification study final report; however, the quantization thresholds used for each of the 14 filter outputs are shown in Figure A-4 and those used for the two regression coefficients c_1 and c_2 are shown below:

5. ENERGY

For each time sample, a measure of the energy was also computed. As an aid to distinguishing vowels from nasals (which usually have most of their energy in a_{1j}) and vowels from sibilants (which usually have most of their energy in a_{14j}), these two filters were not used in computing the energy measure in the following expression:

$$E = \sqrt{\sum_{i=2}^{13} (a_i)^2 - \frac{1}{11} \left(\sum_{i=2}^{13} a_i \right)^2}$$
 (A-7)

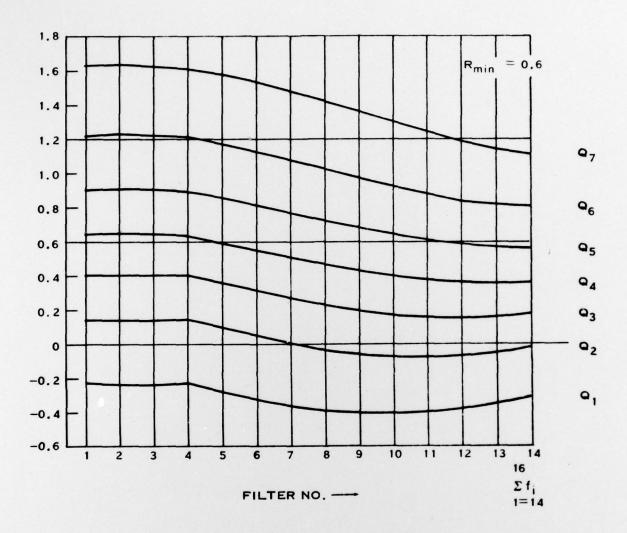


Figure A-4. Quantization Thresholds for Filters

APPENDIX B SCATTER MATRICES

Inherent in many criteria used in clustering is the concept of the scatter of the data, numerically represented by scatter matrices. The within-class scatter matrix measures the distance of the "n" sample vectors from their mean vectors and is the sum of the scatter matrices for all "c" classes. The within-class scatter is given by

$$S_W = \sum_{i=1}^c S_i = \sum_{i=1}^c \sum_{x \in x_i} (\vec{x} - \vec{m}_i) (\vec{x} - \vec{m}_i)^T$$
 (B-1)

The between-class scatter matrix measures the distance of the class mean vectors from the overall sample mean and is given by

$$S_{B} = \sum_{i=1}^{c} n_{i} (\overline{m}_{i} - \overline{m}) (\overline{m}_{i} - \overline{m})^{T}$$
 (B-2)

The total scatter matrix measures the distance of the samples from the overall mean as given by

$$S_{T} = \sum_{\mathbf{x} \in \mathbf{x}} (\mathbf{x} - \mathbf{m}) (\mathbf{x} - \mathbf{m})^{T}$$
 (B-3)

Note that $S_T = S_W + S_B$, and, therefore, $|S_T| = |S_W| + |S_B|$ and tr $S_T = \text{tr } S_W + \text{tr } S_B$.

Anderberg³² summarizes the four principal criteria that have emerged using scatter matrices:

- (1) Minimize tr S_W. This is identical to the sum-of-squared error criterion of Ward.
- (2) Minimize the ratio $|S_W|/|S_T|$. This criterion is known as Wilks' lambda statistic. Equivalent criteria are minimizing $|S_W|$ or maximizing $|S_T|/|S_W|$ or $|I + S_W^{-1}|S_B|$.
- (3) Maximize largest eigenvalue of Sw⁻¹ S_B (attributed to S.N. Roy).
- (4) Maximize tr S_W⁻¹ S_B (Hotelling's trace criterion).

The last three of these criteria involve the eigenvalues of S_W^{-1} S_B that are invariant under nonsingular linear transformations, measuring the ratio of the between-class to within-class scatter in the direction of the eigenvectors. Duda and Hart, however, note that invariant criterion functions are more likely to possess multiple local extrema, and are correspondingly more difficult to extremize.

As noted, the first criterion is simply the sum-of-squared error criterion,

$$J_e = \sum_{i=1}^{c} J_i$$

PRECEDING PAGE NOT FILMED BLANK Specifically,

$$J_{c} = \text{tr } S_{W} = \sum_{i=1}^{c} \text{tr } S_{i} = \sum_{i=1}^{c} \sum_{\vec{x} \in Y_{i}} ||\vec{x} - \vec{m}_{i}||^{2}$$
 (B-4)

With some manipulation, this can be rewritten as

$$J_{e} = \frac{1}{2} \sum_{i=1}^{c} \frac{1}{n_{i}} \sum_{\vec{x} \in X_{i}} \sum_{\vec{x}' \in X_{i}} \|\vec{x} - \vec{x}'\|^{2}$$
(B-5)

Minimization of this criterion is equivalent to maximizing

tr S_B =
$$\sum_{i=1}^{c} n_i ||\vec{m}_i - \vec{m}||^2$$
 (B-6)

since tr S_T is a constant. The term tr S_B can also be manipulated (see Appendix C) to give

tr
$$S_B = \frac{1}{2n} \sum_{i=1}^{c} \sum_{j=1}^{c} n_i n_j ||\vec{m}_i - \vec{m}_j||^2$$
 (B-7)

Although the criteria given above differ, the underlying model using scatter matrices is that of "c" fairly well separated clouds containing roughly equal numbers of points. Duda and Hart, 28 however, demonstrate (Figure B-1) how this model can lead to poor results.

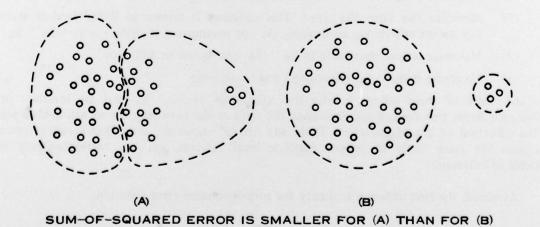


Figure B-1. Problem of Splitting Large Clusters

(FROM DUDA AND HART²⁸)

By definition,

tr
$$S_B = \sum_{i=1}^c n_i ||\vec{m}_i - \vec{m}||^2 = \sum_{i=1}^c n_i (\vec{m}_i - \vec{m})^T (\vec{m}_i - \vec{m})$$
 (C-1)

Expanding,

$$\operatorname{tr} S_{B} = \sum_{i=1}^{c} n_{i} \overline{m}_{i}^{T} \overline{m}_{i} - 2 \sum_{i=1}^{c} n_{i} \overline{m}_{i}^{T} \overline{m} + \sum_{i=1}^{c} n_{i} \overline{m}^{T} \overline{m}$$
 (C-2)

But the sum in the middle term can be rewritten as

$$\sum_{i=1}^{c} n_{i} \vec{m}_{i}^{T} = n \vec{m}^{T} = \sum_{i=1}^{c} n_{i} \vec{m}^{T}$$
 (C-3)

yielding

$$\operatorname{tr} S_{B} = \sum_{i=1}^{c} n_{i} \overline{m}_{i}^{T} \overline{m}_{i} - \sum_{i=1}^{c} n_{i} \overline{m}^{T} \overline{m}$$
 (C-4)

Multiplying by (2n/2n) and splitting the first term in half with appropriate change in indices gives

$$\operatorname{tr} S_{B} = (1/2n) \left[n \sum_{i=1}^{c} n_{i} \vec{m}_{i}^{T} \vec{m}_{i} - 2 \sum_{i=1}^{c} n_{i} m_{i}^{T} \sum_{j=1}^{c} n_{j} \vec{m}_{j} + n \sum_{j=1}^{c} n_{j} \vec{m}_{j}^{T} \vec{m}_{j} \right]$$
 (C-5)

Expressing n as the sum of the n_is (or n_js) and factoring out the summations and n_in_i yields

tr S_B = (1/2n)
$$\sum_{i=1}^{c} \sum_{j=1}^{c} n_i n_j \left[\vec{m}_i^T \vec{m}_i - 2 \vec{m}_i^T \vec{m}_j + \vec{m}_j^T \vec{m}_j \right]$$
 (C-6)

tr S_B =
$$(1/2n)$$
 $\sum_{i=1}^{c} \sum_{j=1}^{c} n_i n_j ||\vec{m}_i - \vec{m}_j||^2$ (C-7)

APPENDIX D

POST ITERATIVE OPTIMIZATION STATISTICS FOR RECOGNITION PATTERNS

STATISTICS FOR DIGIT: 0; REF PT: 0; NO OF DATA PTS: 166
MINAVE AND MINMAX AGGLOM CLUSTERING; 16 MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

C ITE	OF RS	JE (=	TK(W))	JE(C)-JE(C+1 /JE(C)	
AVE 1 75 3 75 4 124 5 134 6 113 7 161 8 136 9 116 10 122	MA X 0 1 3 4 2 0 0 0 5 2 7 5 4 8 5	AVE 69662.7 595867.7 555867.2 551245.6 48248.7 48248.7 44223.7	MAX 69662.7 59345.7 57150.6 54136.8 50045.1 49112.9 48357.3 46051.3*	AVE	8 0.000 0.000 7 0.174 0.174 3 0.247 0.219 0 0.291 0.287 7 0.329 0.341 9 0.359 0.392 7 0.427 0.418 0 0.444 0.444 7 0.467 0.473
C (N-C) T	R(B)	BTL'S		(N-C) *DELJE /JE(C	(N-C) *TR(B)
AVE 1 0.000 2 0.172 3 0.182 4 0.189 5 0.200 6 0.208 7 0.239 8 0.241 9 0.248 10 0.265 STATISTICS MINAVE AND		AGGLOM	CLUSTERI	AVE 24.495 24.495 24.495 24.495 24.689 6.49 3.568 5.95 7.570 2.98 1.864 2.76 2.486 3.11 4.185 4.16 ************************************	6 28.591 28.511 1 20.124 17.843 4 13.702 15.485 4 13.253 13.706 0 11.500 12.544 5 11.313 11.088 5 10.018 10.011 9.163 9.274 8.789 8.887 F DATA PTS: 168 16 MAR79
C NO				15 (C) - 15 (C+1)
C IIE	RS	JE (=	TR(W))	JE(C)-JE(C+1 /JE(C)	TR(B)/TR(W)
AVE 1 0 2 19 3 73 4 86 5 132 6 135 7 155 8 112 9 132 10 166	MA X 0 33 557 754 891 993 83	JE (= 70679.1 60969.0 57069.9 52645.4 50998.3 49329.7 446239.0 45093.1	TR(W)) MAX 70679.1 609197.6 552655.6 50653.9 493522.1 474291.0 44985.8 44201.1	AVE MAX 0.137 0.13 0.064 0.07 0.076 0.06 0.031 0.03 0.033 0.02 0.031 0.03 0.031 0.03 0.031 0.03 0.031 0.03	TR(B)/TR(W) AVE MAX 8 0.000 0.000 7 0.159 0.160 3 0.238 0.258 7 0.343 0.342 6 0.386 0.395 9 0.433 0.432 3 0.482 0.432 8 0.529 0.527
AVE 1 19 2 19 3 73 4 86 5 135 6 135 7 155 8 112	MAX 0 338 557 754 811 933 8 8 (B)	70679.1 60969.9 57069.9 52645.4 50998.3 49329.7 46239.0 45093.1	MAX 70679.1 60914.4 56197.6 52655.6 50683.9 49352.1 47402.8 46291.0 44985.8 44201.1*	AVE MAX 0.137 0.13 0.064 0.07 0.076 0.06 0.031 0.03 0.033 0.02 0.031 0.03 0.031 0.03 0.031 0.03 0.031 0.03	TR(B)/TR(W) AVE MAX 0.000 0.000 0.159 0.160 0.238 0.258 0.343 0.342 6 0.386 0.395 9 0.433 0.432 0.432 0.432 0.529 0.527 0.567 0.571 * 0.604 0.599 (N-C)*TR(B) /(C-1)*TR(W)

STATISTICS FOR DIGIT: 2; REF PT: 0; NO OF DATA PTS: 166
MINAVE AND MINMAX AGGLOM CLUSTERING; 16 MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO OF	JE (=TR(W))	\\TE(C) -\TE(C)	TR(B)/TR(W)
12345678910	AVE MAX 0 0 72 48 105 48 140 50 132 82 92 124 116 121 134 133 102 113 138 83	AVE MAX 51525.5 51525.5 45108.0 45140.0 42371.1 42424.6 40377.6 41131.9 39269.9 39098.7 37929.8 37788.9 36991.1 36777.1 35861.8 35640.8 35323.3 35122.3 34182.5 34613.7	AVE MAX 0.125 0.124 0.061 0.060 0.047 0.030 0.027 0.049 0.034 0.033 0.025 0.027 0.031 0.031 0.015 0.015 0.032 0.014	AVE 0.000 0.000 0.142 0.141 0.216 0.215 0.276 0.253 0.312 0.318 0.358 0.368 0.393 0.401 0.437 0.446 0.459 0.467 0.507 0.489
С	C(N-C)TR(B) /2N(C-1)TR(W)	HTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) * [R(B) /(C-1) * [R(A)
12345678910	AVE MAX 0.000 0.000 0.140 0.139 0.158 0.157 0.179 0.163 0.188 0.192 0.206 0.209 0.218 0.223 0.236 0.241 0.243 0.247 0.263 0.254	AVE MAX 0.000 0.000 0.285 0.284 0.328 0.329 0.363 0.437 0.385 0.486 0.496 0.494 0.545 0.499 0.548 0.539 0.608 0.626	AVE MAX 20.675 20.572 10.011 9.926 7.716 4.997 4.472 8.057 5.528 5.427 3.985 4.311 4.885 4.943 2.388 2.313 5.102 2.288	AVE 0.000 23.474 17.716 17.591 15.001 12.640 12.640 12.872 11.542 11.705 10.478 10.643 9.921 9.921 9.921 9.921 9.921 9.921 9.921 9.921 9.921 9.921 9.923 9.859 8.851

STATISTICS FOR DIGIT: 3; REF PT: 0; NO OF DATA PTS: 169
MINAVE AND MINMAX AGGLOM CLUSTERING; 16 MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO OF ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)
1234567890	AVE MAX 0 0 70 51 98 70 113 49 118 67 90 68 76 66 106 58 112 82 127 67	AVE 54075.5 47673.7 47698.1 44988.2 45016.6 43089.4 43089.4 41295.5 39523.6 39708.2 38841.5 38845.5 36810.8 36009.3 35675.6 35267.3	AVE MAX 0.118 0.118 0.056 0.056 0.042 0.045 0.050 0.040 0.035 0.038 0.017 0.042 0.033 0.031 0.020 0.024 0.031 0.021	AVE
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C) * TR(B) /(C-1) * TR(A)
12345678910	AVE 0.000 0.000 0.133 0.132 0.149 0.148 0.166 0.168 0.213 0.209 0.219 0.236 0.239 0.254 0.250 0.267 0.270 0.279	AVE MAX 0.000 0.000 0.270 0.270 0.302 0.295 0.365 0.348 0.392 0.382 0.458 0.495 0.493 0.490 0.566 0.547 0.594 0.550 0.639 0.579	AVE 19.889 19.813 9.407 9.386 7.006 7.410 8.196 6.567 5.709 6.304 2.813 6.834 5.322 4.957 3.234 3.297	AVE 0.000 0.000 22.425 22.328 15.765 16.702 14.023 14.155 13.143 12.688 12.003 11.795 10.590 11.378 10.108 10.725 9.380 10.034 9.112 9.422

STATISTICS FOR DIGIT: 4; REF PT: 0; NO OF DATA PTS: 167
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO UF ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(W)
1234567890	AVE MAX 0 0 99 32 84 64 133 59 124 123 113 107 128 116 141 83 149 129 200 115	AVE MAX 43604.4 43604.4 38061.0 38075.8 35993.5 36052.4 34536.4 34738.5 33525.7 33340.8 32389.8 32054.9 31074.0 31056.4 30165.9 29229.9 29153.2 28899.5	0.015 0.025 0.015 0.035 0.034 0.011	AVE MAX 0.000 0.000 0.146 0.145 0.211 0.209 0.263 0.255 0.301 0.308 0.346 0.360 0.403 0.404 0.424 0.440 0.445 0.492 0.496 0.509
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C) *TR(B) /(C-1) *TR(W)
12345678910	AVE MAX 0.000 0.000 0.140 0.139 0.151 0.150 0.166 0.161 0.177 0.181 0.194 0.202 0.218 0.219 0.224 0.232 0.230 0.253 0.251 0.257	AVE	AVE	AVE 0.000 0.000 23.303 23.232 16.810 16.653 13.828 13.442 11.800 12.083 10.803 11.241 10.417 10.437 9.331 9.671 8.520 9.405 8.372 8.594

STATISTICS FOR DIGIT: 5; REF PT: 0; NO OF DATA PTS: 169
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO OF ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TK(B)/TK(W)
123 45 67 89 10	AVE MAX 0 15 19 52 56 85 45 108 46 127 43 126 59 150 68 167 75 178 84	AVE MAX 87861.3 87861.3 73448.7 73443.4 70442.7 69540.9 66955.3 66296.3 63947.1 64185.0 62144.9 62522.0 60578.4 60714.5 58523.0 59195.0 57107.8 57430.3 55926.8 56100.6	AVE MAX 0.164 0.164 0.041 0.053 0.050 0.047 0.045 0.032 0.028 0.026 0.025 0.026 0.025 0.025 0.024 0.030 0.021 0.023	AVE MAX 0.000 0.000 0.196 0.196 0.247 0.263 0.312 0.325 0.374 0.369 0.414 0.405 0.450 0.444 0.539 0.530 0.571 0.566
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) *TR(B) /(C-1) *TR(W)
1234567890	AVE MAX 0.000 0.000 0.193 0.193 0.181 0.193 0.202 0.210 0.225 0.223 0.238 0.233 0.250 0.248 0.271 0.262 0.285 0.280 0.297 0.294	AVE MAX 0.000 0.000 0.390 0.390 0.335 0.368 0.383 0.403 0.441 0.414 0.467 0.438 0.509 0.542 0.561 0.577 0.611 0.597	AVE 27.404 6.794 8.820 8.169 7.699 7.368 5.223 4.594 4.223 4.084 4.683 5.463 4.770 3.288 3.661	AVE 0.000 32.574 32.588 20.400 17.782 17.069 17.782 13.408 13.131 12.085 11.969 10.703 10.531 10.024

STATISTICS FOR DIGIT: 6; REF PT: 0; NO OF DATA PTS: 167
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE DETIMIZATION FOR MIN JE

C	NO OF ITERS	JE (=TR(A))	JE(C)-JE(C+1) /JE(C)	[k(B)/[k(W)
12345678910	AVE 0 0 84 47 84 27 68 42 88 42 100 47 133 49 130 90 160 32	AVE MAX 58412.0 58412.0 48111.0 48136.5 45669.2 45686.7 44419.0 43177.7 41809.1 42246.4 40753.9 40986.8 39326.5 40228.1 38659.4 38769.0 37594.0 38654.9 36996.0 37826.5	0.027 0.055 0.059 0.022 0.025 0.030 0.035 0.019 0.017 0.036 0.028 0.003 0.016 0.021	AVE 0.000 0.000 0.214 0.213 0.279 0.279 0.315 0.353 0.397 0.363 0.433 0.425 0.485 0.452 0.511 0.506 0.554 0.511 0.579 0.544
С	C(N-C) TR(B) /2N(C-1) TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C) * [R(B) /(C-1) * [R(W)
12345 6789 10	AVE MAX 0.000 0.000 0.209 0.208 0.203 0.203 0.203 0.229 0.238 0.229 0.248 0.243 0.268 0.249 0.274 0.272 0.291 0.269 0.298 0.281	AVE MAX 0.000 0.000 0.422 0.422 0.430 0.426 0.371 0.518 0.475 0.530 0.519 0.486 0.565 0.546 0.572 0.601 0.590 0.601 0.608 0.633	AVE	AVE 0.000 0.000 54.900 34.795 22.601 22.561 16.906 18.935 15.884 15.306 13.779 13.520 12.780 11.905 11.460 11.346 10.798 9.967 9.969 9.372

STATISTICS FOR DIGIT: 7; REF PT: 0; NO DF DATA PTS: 168
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

C	NO OF ITERS	JE (=[H(W))	JE(C)-JE(C+1) /JE(C)	[K(B)/[K(W)
123 45 67 89 10	AVE MAX 0 0 81 2 32 17 71 36 60 67 85 72 93 62 106 79 104 45 98 43	AVE MAX 80124.8 80124.8 64595.8 64656.4 60652.0 60524.4 58433.6 58724.0 57262.1 55331.3 54819.4 53704.3 53809.7 52160.0 51550.4 50680.3 49846.8 50085.5 48808.2 49334.1	AVE	AVE MAX 0.000 0.000 0.240 0.239 0.321 0.324 0.371 0.364 0.399 0.448 0.462 0.492 0.489 0.536 0.554 0.561 0.607 0.600 0.642 0.624
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C) *DELJE /JE(C)	(N-C)*TR(B) /(C-1)*TR(A)
12345 6789 10	AVE 0.000 0.000 0.000 0.236 0.235 0.235 0.237 0.240 0.237 0.241 0.270 0.265 0.283 0.272 0.298 0.300 0.314 0.321 0.317 0.333 0.324	AVE MAX 0.000 0.000 0.475 0.475 0.475 0.485 0.532 0.545 0.568 0.659 0.655 0.671 0.727 0.672 0.717 0.814	AVE MAX 32.172 32.047 10.074 10.545 5.998 4.878 3.268 9.417 6.911 4.763 2.965 4.630 6.718 4.539 5.254 1.860 3.292 2.370	AVE 0.000 59.666 39.474 26.327 26.555 20.169 19.801 16.170 18.148 14.864 15.641 13.041 14.297 12.591 13.197 11.997 11.645 11.193 10.898

STATISTICS FOR DIGIT: B; REF PT: U; NO OF DATA PTS: 168
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE DETIMIZATION FOR MIN JE

C	NO OF ITERS	JE (=TK(W))	JE(C)-JE(C+1) /JE(C)	TR(B)/TR(A)
1234567890	AVE MAX 90 13 135 26 109 50 100 63 102 55 133 48 123 52 140 44 123 58	AVE MAX 61486.2 61486.2 51255.8 51286.6 47621.8 47656.7 43789.9 43799.9 42426.6 42403.3 41395.2 41195.1 40203.8 40101.5 39474.0 38918.9 37863.9 38054.1 36747.5 36408.3	0.071 0.071 0.080 0.081 0.031 0.032 0.024 0.028 0.029 0.027 0.018 0.029 0.041 0.022 0.029 0.043	AVE MAX 0.000 0.000 0.200 0.199 0.291 0.290 0.404 0.404 0.449 0.450 0.485 0.493 0.529 0.538 0.529 0.538 0.624 0.616 0.673 0.689
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) *TR(B) /(C-1) *TR(W)
1 2 3 4 5 6 7 8 9 10	AVE MAX 0.000 0.000 0.197 0.197 0.214 0.214 0.263 0.263 0.272 0.273 0.281 0.285 0.296 0.298 0.303 0.316 0.332 0.328 0.352 0.360	AVE MAX 0.000 0.000 0.394 0.394 0.429 0.429 0.528 0.528 0.571 0.571 0.571 0.579 0.657 0.628 0.628 0.669 0.675 0.743 0.734 0.805	AVE MAX 27.786 27.703 11.769 11.749 13.277 13.353 5.106 5.229 3.963 4.644 4.662 4.301 2.922 4.748 6.527 3.555 4.668 6.676 **************	AVE 0.000 0.000 33.133 33.013 24.019 23.941 22.092 22.074 18.306 18.359 14.205 14.309 12.746 13.254 12.400 12.238 11.819 12.092

STATISTICS FOR DIGIT: 9; REF PT: 0; NO OF DATA PTS: 169
MINAVE AND MINMAX AGGLOM CLUSTERING; 20MAR79
POST ITERATIVE OPTIMIZATION FOR MIN JE

С	NO OF ITERS	JE (=TR(W))	JE(C)-JE(C+1) /JE(C)	TR(b)/TR(w)
12345678910	AVE MAX 0 0 74 15 81 27 128 46 150 110 140 62 117 110 121 58 107 70 141 76	AVE 75929.8 75929.8 65731.1 65777.5 61771.8 61677.0 58517.1 58556.2 55207.2 55211.1 52482.7 53535.6 51106.2 51086.3 49760.9 50559.5 47428.2 48322.0 46395.2 47105.9	AVE MAX 0.134 0.134 0.060 0.062 0.053 0.051 0.057 0.057 0.049 0.030 0.026 0.046 0.026 0.040 0.026 0.045 0.025 0.025	AVE MAX 0.000 0.000 0.155 0.154 0.229 0.231 0.298 0.297 0.375 0.375 0.447 0.418 0.486 0.486 0.526 0.571 0.637 0.612
С	C(N-C)TR(B) /2N(C-1)TR(W)	BTL'S SIGMA	(N-C)*DELJE /JE(C)	(N-C) * TR(B) /(C-1) * TR(A)
12345678910	AVE MAX 0.000 0.000 0.153 0.153 0.169 0.170 0.194 0.193 0.228 0.228 0.259 0.242 0.272 0.272 0.286 0.273 0.320 0.304 0.333 0.320	AVE 0.000 0.000 0.307 0.307 0.345 0.342 0.392 0.392 0.477 0.458 0.536 0.497 0.549 0.582 0.611 0.597 0.683 0.657 0.721 0.728	AVE 22.463 10.059 10.411 3.746 8.399 9.333 9.426 8.093 4.977 4.275 7.457 4.264 1.671 7.547 7.125 3.485 4.026	AVE MAX 0.000 0.000 25.911 25.775 19.023 19.180 16.366 16.318 15.390 15.38 14.564 13.637 13.115 13.130 12.096 11.541 12.019 11.427 11.246 10.810

MISSION of Rome Air Development Center

actures content conten

RADC plans and executes research, development, test and selected acquisition programs in support of Command, Control Communications and Intelligence (C³I) activities. Technical and engineering support within areas of technical competence is provided to ESD Program Offices (POs) and other ESD elements. The principal technical mission areas are communications, electromagnetic guidance and control, surveillance of ground and aerospace objects, intelligence data collection and handling, information system technology, ionospheric propagation, solid state sciences, microwave physics and electronic reliability, maintainability and compatibility.